

# GigaScience

## Population genomics of wild Chinese rhesus macaques reveals a dynamic demographic history and local adaptation, with implications for biomedical research --Manuscript Draft--

<b>Manuscript Number:</b>	GIGA-D-17-00291R2													
<b>Full Title:</b>	Population genomics of wild Chinese rhesus macaques reveals a dynamic demographic history and local adaptation, with implications for biomedical research													
<b>Article Type:</b>	Research													
<b>Funding Information:</b>	<table border="1"><tr><td>Key Project of National Natural Science Foundation of China (31530068)</td><td>Dr Ming Li</td></tr><tr><td>Strategic Priority Research Program of the Chinese Academy of Sciences (XDA19050202)</td><td>Dr Ming Li</td></tr><tr><td>Strategic Priority Research Program of the Chinese Academy of Sciences (XDB31000000)</td><td>Dr Ming Li</td></tr><tr><td>National Key R&amp;D Program of China (2016YFC0503200)</td><td>Dr Ming Li</td></tr><tr><td>Creative Research Group Project of NSFC; and Science &amp; Technology Department of Sichuan Province (2018JZ0008)</td><td>Dr Ming Li</td></tr><tr><td>Natural Science Foundation of China (31471989)</td><td>Dr Zhijin Liu</td></tr></table>		Key Project of National Natural Science Foundation of China (31530068)	Dr Ming Li	Strategic Priority Research Program of the Chinese Academy of Sciences (XDA19050202)	Dr Ming Li	Strategic Priority Research Program of the Chinese Academy of Sciences (XDB31000000)	Dr Ming Li	National Key R&D Program of China (2016YFC0503200)	Dr Ming Li	Creative Research Group Project of NSFC; and Science & Technology Department of Sichuan Province (2018JZ0008)	Dr Ming Li	Natural Science Foundation of China (31471989)	Dr Zhijin Liu
Key Project of National Natural Science Foundation of China (31530068)	Dr Ming Li													
Strategic Priority Research Program of the Chinese Academy of Sciences (XDA19050202)	Dr Ming Li													
Strategic Priority Research Program of the Chinese Academy of Sciences (XDB31000000)	Dr Ming Li													
National Key R&D Program of China (2016YFC0503200)	Dr Ming Li													
Creative Research Group Project of NSFC; and Science & Technology Department of Sichuan Province (2018JZ0008)	Dr Ming Li													
Natural Science Foundation of China (31471989)	Dr Zhijin Liu													
<b>Abstract:</b>	<p><b>Background:</b> The rhesus macaque (RM, <i>Macaca mulatta</i>) is the most important nonhuman primate model in biomedical research. We present the first genomic survey of wild RMs, sequencing 81 geo-referenced individuals of five subspecies from 17 locations in China, a large fraction of the species' natural distribution.</p> <p><b>Results:</b> Populations were structured into five genetic lineages on the mainland and Hainan Island, recapitulating current subspecies designations. These subspecies are estimated to have diverged 125.8 to 51.3 thousand years ago, but feature recent gene flow. Consistent with the expectation of a larger body size in colder climates and smaller body size in warmer climates (Bergman's rule), the northernmost RM lineage (<i>M. m. tcheliensis</i>), possessing the largest body size of all Chinese RMs, and the southernmost lineage (<i>M. m. breviceaudus</i>), with the smallest body size of all Chinese RMs, feature positively selected genes responsible for skeletal development. Further, two candidate selected genes (<i>Fbp1</i>, <i>Fbp2</i>) found in <i>M. m. tcheliensis</i> are involved in gluconeogenesis, potentially maintaining stable blood glucose levels during starvation when food resources are scarce in winter. The tropical subspecies <i>M. m. breviceaudus</i> showed positively selected genes related to cardiovascular function and response to temperature stimuli, potentially involved in tropical adaptation. We found 118 SNPs matching human disease-causing variants with 82 being subspecies-specific.</p> <p><b>Conclusions:</b> These data provide a resource for selection of RMs in biomedical experiments. The demographic history of Chinese RMs, and their history of local adaptation offers new insights into their evolution and provides valuable baseline information for biomedical investigation.</p>													
<b>Corresponding Author:</b>	Ming Li  CHINA													
<b>Corresponding Author Secondary Information:</b>														
<b>Corresponding Author's Institution:</b>														
<b>Corresponding Author's Secondary Institution:</b>														

<b>First Author:</b>	Zhijin Liu
<b>First Author Secondary Information:</b>	
<b>Order of Authors:</b>	Zhijin Liu
	Xinxin Tan
	Pablo Orozco-terWengel
	Xuming Zhou
	Liye Zhang
	Shilin Tian
	Zhongze Yan
	Huailiang Xu
	Baoping Ren
	Peng Zhang
	Zuofu Xiang
	Binghua Sun
	Christian Roos
	Michael W. Bruford
	Ming Li
<b>Order of Authors Secondary Information:</b>	
<b>Response to Reviewers:</b>	<p>Dear editor,</p> <p>We thank you for your encouragement and help to our manuscript “Population genomics of wild Chinese rhesus macaques reveals a dynamic demographic history and local adaptation, with implications for biomedical research”.</p> <p>We have shorted the abstract and structured it in sections (“Background – Results - Conclusions”).</p> <p>We have added a citation to our upcoming GigaDB dataset without the DOI at the end of the reference list and cited this in the ‘Data Access’.</p> <p>The point-to-point responses to the reviewers are listed below.</p> <p>Reviewer #1:  All my major concerns have been satisfactorily addressed. However, I see several typos or other minor problems with the text.</p> <p>Comment_1-1:  line 162: should be "from" instead of "form"  line 230: I think the authors mean "upregulated" not "unregulated"  line 285: should be "have" not "has"  I think there are other similar issues.</p> <p>Response_1-1:  These issues have been amended accordingly. (line 158, line 226, line 283)</p> <p>Reviewer #2:  Comment_2-1:  The authors have revised their manuscript based on the detailed reviewer reports and have substantially improved it in the process. Most of my previous methodological concerns have been addressed. However, there still seems to be some confusion regarding my comments about proper calculation of genetic diversity (comment 2-3</p>

and 2-4). The authors write in their rebuttal that they have used the GATK Best Practices workflow with multi-sample genotyping, thereby solving the problem of uncallable sites by accumulating evidence across samples. This is only partly mitigating the problem here. If a site is not callable across all samples due to poor mappability or a 'N' in the reference sequence, a genetic variant cannot be detected at this position. Thus, genetic diversity will be underestimated when assuming that variant sites can occur along the entire genome. To obtain an accurate estimate of genetic diversity, it is therefore crucial to consider the exact number of sites that are callable, i.e. where it is possible to detect a potential variant. This can be achieved for example with GATK's 'CallableLoci' module. Alternatively, the '--includeNonVariantSites' flag of GATK's 'GenotypeGVCFs' can be used to emit both confident variant and non-variant sites. To be considered in the calculation of genetic diversity, a site should be callable (i.e. have a valid variant or non-variant genotype) in a certain proportion (e.g. 80%) of individuals. This initial site filter has to be independent of the variant state of a site, i.e. also variant sites in the SNP data set have to be filtered out if they don't fulfill the callability criteria.

Response\_2-1:

We followed this helpful suggestion and redone the 'GenotypeGVCFs' in GATK with the '--includeNonVariantSites' flag to get both the variant and non-variant sites. Besides the basic hard filter by 'VariantFiltration' in GATK, we also filtered out the variants with a 'N' in the reference sequence or the sites including more than 20% missing genotypes. For the non-variant sites, we did the same filter and retain only the callable non-variant sites. The genetic diversity and heterozygosity have been re-estimated based on all the callable sites. (lines 101-102, lines 109-112, lines 379-390 and Table 1)

Comment\_2-2:

Lines 76-79: These sentences are unclear. If to date only 9 captive Chinese RMs have been sequenced, how could Zhong et al. assess genetic diversity in 26 Chinese individuals? I guess the authors mean "Until recently, ..." rather than "To date, ..." at the beginning of the first sentence.

Response\_2-2:

This issue has been amended accordingly. (lines 74-77)

Comment\_2-3:

Line 102: This sentence is confusing. The authors write that they identified ~58 mio SNPs in the 81 Chinese RMs. From their explanations, I understood that this is the total number of variant sites, i.e. including fixed differences to the reference genome? The reference genome is of Indian origin, so it's incorrect to write that these are SNPs in Chinese RMs.

Response\_2-3:

This issue has been amended accordingly. We have filtered out the fixed differences to the reference genome. (lines 386-388)

Comment\_2-4:

Line 103: Was Watterson's theta correctly estimated considering only sites actually segregating within the Chinese RMs?

Response\_2-4:

The  $\theta_W$  and  $\theta_{\pi}$  have been re-estimated only based on segregating variations within Chinese RMs. (lines 101-102 and Table 1)

Comment\_2-5:

Line 104: "and the nucleotide diversity measured by segregating sites (Watterson's  $\theta$ ,  $\theta_W$ ) and mean pairwise differences ( $\theta_{\pi}$ ) is ..."

Response\_2-5:

This issue has been amended accordingly. (lines 101-102)

Comment\_2-6:

Lines 106-110: It doesn't make sense to use variant sites relative to a reference genome for the analysis of shared and private SNPs. These numbers reflect a mixture of segregating variation and fixed differences to the reference genome. Please redo these analyses by only considering actual segregating variation within the compared entities.

Response\_2-6:

We have filtered out the fixed difference to the reference genome. All these analyses have redone based on actual segregating variations within Chinese RMs. (lines 104-108, lines 387-388)

Comment\_2-7:

Line 139: "based on  $\theta W$  and  $\theta\pi$  are ..."

Response\_2-7:

This issue has been amended accordingly. (line 137)

Comment\_2-8:

Lines 170-177: Round estimates and provide confidence intervals.

Response\_2-8:

This issue has been amended accordingly. (lines 166-173)

Comment\_2-9:

Lines 180-181: Tone down this statement, since you haven't explicitly compared models with and without gene flow. Something along the lines of: "Our results indicate that low levels of gene flow occurred between all five extant lineages of Chinese RMs."

Response\_2-9:

This issue has been amended accordingly. Substantial gene flows have been detected between different subspecies. Please see the response to the comment\_2-20 and Supplementary Table 5. (lines 177-178)

Comment\_2-10:

Line 193: "led to further differentiation by limiting gene flow among them."

Response\_2-10:

This issue has been amended accordingly. (line 190)

Comment\_2-11:

Lines 208-211: This sentence seems to conflict with the sentence on lines 200-204.

Response\_2-11:

For *M. m. tcheliensis*, which occurs in the northernmost range of the RMs under cold conditions, we first estimated  $F_{ST}$  and  $\theta\pi$  between it and each of the other four subspecies. Then we got four lists of candidate genes in *M. m. tcheliensis*. The final positive selection genes are the intersection of these four lists. Similar, in the case of *M. m. brevicaudus*, we used the same method. The details of this process are shown in Supplementary Fig. 7. We chose this method to get the final positive selection genes, instead of directly comparing *M. m. tcheliensis* and *M. m. brevicaudus*, for the purpose of reducing the false positives of the results and obtaining more accurate selective gene lists. (lines 197-201, lines 205-208 and Supplementary Fig. 8)

Comment\_2-12:

Lines 223-226: See previous comment 2-16.

Response\_2-12:

This issue has been amended accordingly. (lines 220-223)

Comment\_2-13:

Line 230: "upregulated" instead of "unregulated"?

Response\_2-13:

This issue has been amended accordingly. (line 226)

Comment\_2-14:

Lines 240-241: Having long forearms doesn't really fit the expectation, as long extremities would increase the surface to volume ration. I realize that forearm length is probably strongly correlated with body size, but this is confusing for the reader. Maybe just omit the forearm length.

Response\_2-14:

This issue has been amended accordingly. We have omitted the forearm length in the revised manuscript. Many thanks for this helpful suggestion. (lines 237-238)

Comment\_2-15:

Line 385-387: Provide details about the filter settings. The current description of the variant hard filtering approach doesn't allow to reproduce the data set used for the downstream analyses.

Response\_2-15:

This issue has been amended accordingly. After variant calling, we first applied the “SelectVariants” to exclude the Indel and split the variant and non-variant sites. Then we applied the hard filter command ‘VariantFiltration’ to exclude potential false-positive variant calls with the following criteria: “-filterExpression ‘QD < 5.0 || FS > 60.0 || MQ < 40.0 || ReadPosRankSum < -8.0 || MQRankSum < -12.5’” and “--genotypeFilterExpression ‘DP < 4.0’”. Additionally, the sites are filtered if there is a ‘N’ is in the reference sequence; if the site is fixed difference to the reference genome or if the site including more than 20% missing genotypes. (lines 383-388)

Comment\_2-16:

Line 410-411: Provide details of how the consensus sequences have been generated.

Response\_2-16:

This issue has been amended accordingly. We called the consensus sequences using Samtools mpileup [68] by applying: “samtools mpileup -q 1 -C 50 -S -D -m 2 -F 0.002 -u -f \*.fa(genome) \*.bam | bcftools view -c - | vcfutils.pl vcf2fq -d 10 -D 100 -Q 20 -> \*.psmc.fq” and “fq2psmcfa -q10 -s 100 \*.psmc.fq >\*.psmc.fa”. To ensure the quality of consensus sequences, we used data of ten individuals with an average coverage >20x (22.20-34.32x). (lines 411-415)

Comment\_2-17:

Line 418: Provide more details of how the SNP data has been converted to joint site frequency spectra. How was the number of non-variant sites assessed accurately (see comment above)?

Response\_2-17:

VCF file containing callable variant sites was used converted to fastsimcoal style folded SFS. To mitigate the effect of linkage disequilibrium, we filtered out the SNPs located within 10 kb from genes and then we took one SNPs every 10kb randomly. The multidimensional folded SFS for all the five subspecies is generated by easySFS (<https://github.com/isaacovercast/easySFS#easysfs>). The non-variant sites were not used to convert the SFS. (lines 426-427)

Comment\_2-18:

Lines 422-423: Not sure what this is supposed to mean. Have you simulated data sets under the inferred model and compared distributions of simulated summary statistics to the observed values? However, Supplementary Table 5 doesn't show distributions of summary statistics, rather estimates of model parameters. Provide details of how confidence intervals have been calculated and show how good the model fits the observed data.

Response\_2-18:

Lines 422-423 is a typo-error and has been removed. We got confidence intervals after parameter estimation using parametric bootstraps. We chose the replicate with the highest estimated maximum likelihood to generate parametric bootstraps. One hundred multidimensional SFS files were generated for this set of parameters and then estimated parameters from these pseudo-observed data sets using the same tpl and est files as those used to get the parameters with highest likelihood. We used the option ‘-initvalues file.pv’ to reduce the number of runs necessary to estimate parameters when estimating confidence intervals by bootstrap. The ‘file.pv’ containing initial parameter values for parameter estimation is automatically generated after parameter estimation by fsc26. The observed data and the confidence intervals from 100 parametric bootstraps were showed in Supplementary Fig. 7 and Supplementary Table 5. (lines 427-437)

Comment\_2-19:

Line 475: I haven't been able to find any notes in the Supplementary Information.

Response\_2-19:

This issue has been amended accordingly. According to the format of GigaScience, the supplementary files does not include any notes. This confusion on line 475 is a typo and we have modified it (line 491).

Comment\_2-20:

Supplementary Table 5: Are the gene flow estimates really representing the number of migrants (i.e. Nem)? This would be completely negligible gene flow. Or are these

	<p>numbers rather migration rates (i.e. m), which would imply quite substantial gene flow. Response_2-20:  It is a typo-error. The gene flow estimated really represent the migration rate between subspecies, which imply quite substantial gene flow (Supplementary Table 5).</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
<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 deposited in <a href="#">publicly available repositories</a> (where available and ethically appropriate), referencing such data using</p>	Yes

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 **Title:** Population genomics of wild Chinese rhesus macaques reveals a  
2 dynamic demographic history and local adaptation, with implications for  
3 biomedical research

4 **Running Title:** Population genomics of wild rhesus macaques

5 Zhijin Liu<sup>1, †</sup>, Xinxin Tan<sup>1, 2, †</sup>, Pablo Orozco-terWengel<sup>3, †</sup>, Xuming Zhou<sup>1, 4</sup>, Liye Zhang<sup>1, 2</sup>,

6 Shilin Tian<sup>5</sup>, Zhongze Yan<sup>1, 6</sup>, Huailiang Xu<sup>7</sup>, Baoping Ren<sup>1</sup>, Peng Zhang<sup>8</sup>, Zuofu Xiang<sup>9</sup>,

7 Binghua Sun<sup>10</sup>, Christian Roos<sup>11</sup>, Michael W. Bruford<sup>3, \*</sup>, Ming Li<sup>1, 12 \*</sup>

8 <sup>1</sup> CAS Key Laboratory of Animal Ecology and Conservation Biology, Institute of Zoology,  
9 Beijing, China.

10 <sup>2</sup> University of Chinese Academy of Sciences, Beijing 100039, China.

11 <sup>3</sup> School of Biosciences, Cardiff University, Sir Martin Evans Building, Museum Avenue, Cardiff  
12 CF10 3AX, United Kingdom.

13 <sup>4</sup> Division of Genetics, Department of Medicine, Brigham and Women's Hospital, Harvard  
14 Medical School, Boston, MA 02115, USA.

15 <sup>5</sup> Novogene Bioinformatics Institute, Beijing 100083, China.

16 <sup>6</sup> Institute of Physical Science and Information Technology, Anhui University, Hefei, 230601,  
17 China.

18 <sup>7</sup> College of Life Science, Sichuan Agricultural University, Ya'an 625014, China.

19 <sup>8</sup> School of Sociology and Anthropology, Sun Yat-sen University, Guang Zhou, China.

20 <sup>9</sup> College of Life Science and Technology, Central South University of Forestry and Technology,  
21 Changsha 410004, Hunan, China.

22 <sup>10</sup> School of Life Sciences, Anhui University, Hefei, 230601, China.

23 <sup>11</sup> Gene Bank of Primates and Primate Genetics Laboratory, German Primate Center, Leibniz  
24 Institute for Primate Research, Kellnerweg 4, 37077 Göttingen, Germany.

25 <sup>12</sup> Center for Excellence in Animal Evolution and Genetics, Chinese Academy of Sciences,  
26 Kunming, 650223, China.

27 <sup>†</sup> Contributed equally

28 \* Correspondence: Ming Li, [lim@ioz.ac.cn](mailto:lim@ioz.ac.cn); Michael W. Bruford, [BrufordMW@cardiff.ac.uk](mailto:BrufordMW@cardiff.ac.uk)



1      29    **Abstract**

2  
3      30    **Background:** The rhesus macaque (RM, *Macaca mulatta*) is the most important nonhuman  
4  
5      31    primate model in biomedical research. We present the first genomic survey of wild RMs,  
6  
7      32    sequencing 81 geo-referenced individuals of five subspecies from 17 locations in China, a large  
8  
9      33    fraction of the species' natural distribution.

10  
11     34    **Results:** Populations were structured into five genetic lineages on the mainland and Hainan  
12  
13     35    Island, recapitulating current subspecies designations. These subspecies are estimated to have  
14  
15     36    diverged 125.8 to 51.3 thousand years ago, but feature recent gene flow. Consistent with the  
16  
17     37    expectation of a larger body size in colder climates and smaller body size in warmer climates  
18  
19     38    (Bergman's rule), the northernmost RM lineage (*M. m. tcheliensis*), possessing the largest body  
20  
21     39    size of all Chinese RMs, and the southernmost lineage (*M. m. breviceaudus*), with the smallest  
22  
23     40    body size of all Chinese RMs, feature positively selected genes responsible for skeletal  
24  
25     41    development. Further, two candidate selected genes (*Fbp1*, *Fbp2*) found in *M. m. tcheliensis* are  
26  
27     42    involved in gluconeogenesis, potentially maintaining stable blood glucose levels during  
28  
29     43    starvation when food resources are scarce in winter. The tropical subspecies *M. m. breviceaudus*  
30  
31     44    showed positively selected genes related to cardiovascular function and response to temperature  
32  
33     45    stimuli, potentially involved in tropical adaptation. We found 118 SNPs matching human  
34  
35     46    disease-causing variants with 82 being subspecies-specific.

36  
37  
38     47    **Conclusions:** These data provide a resource for selection of RMs in biomedical experiments.  
39  
40     48    The demographic history of Chinese RMs, and their history of local adaption offers new insights  
41  
42     49    into their evolution and provides valuable baseline information for biomedical investigation.

43  
44     50    **Keywords:** *Macaca mulatta*, population genomics, adaptive selection, biomedical model  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

## 51 **Introduction**

52 Understanding how species evolve and adapt to their environments is an essential question in  
53 evolutionary biology. Rhesus macaques (RMs, *Macaca mulatta*) are, after humans, the world's  
54 most widely distributed primates [1-5], occupying a vast geographic distribution spanning from  
55 Afghanistan to the Chinese shore of the Pacific Ocean and south into Myanmar, Thailand, Laos  
56 and Vietnam [5]. As the most widely distributed nonhuman primate species, RMs occupy diverse  
57 ecological landscapes and habitats, making them an interesting model to address questions about  
58 how species evolve and adapt to local environmental variation, including characterizing the  
59 genomic architecture of adaptation to habitat, climate and other biotic and abiotic factors. Yet,  
60 despite much work on primate comparative genomics, very few population genomic studies have  
61 been carried out on wild RMs [6, 7]. Importantly, as RMs are widely used as a primate model in  
62 physiological, psychological and cognitive studies [8-10], knowledge about their genomic  
63 architecture could improve and refine biomedical research [10] as the genomic composition of  
64 experimental animals can have a considerable influence on the outcome of experiments [11, 12].  
65 Therefore, information on the genomic diversity not only of captive, but also of wild RMs, that  
66 could become a genomic resource for future utilization in medical research, is essential.

67 In biomedical research, two main RM populations (Indian and Chinese) are recognized [6,  
68 13]. They diverged from each other ~162 thousand years ago (kya) and are characterized by  
69 extensive differences in morphology, behavior, ecology, physiology, reproduction, and disease  
70 progression [6, 13-19]. In 1978 India banned all RM exports to breeding centers across the world,  
71 thus curtailing the availability of wild Indian RMs and subsequently increasing the demand for  
72 Chinese RMs in biomedical research, thereby making a detailed characterization of genetic  
73 variants from Chinese RMs crucial for biomedical usage of this species.

74 Until recently, the genomes of 133 captive RMs from eight colonies have been sequenced,  
75 however, 124 of them are of Indian-origin and only nine individuals were presumed to be of  
76 Chinese origin [6]. Besides, Zhong *et al.* [7] reported genomic variation in 26 Chinese captive  
77 RMs identifying ~46 million (M) single nucleotide polymorphisms (SNPs). Nevertheless, most  
78 of the RM genetic variation known to date is limited to captive populations which may contain  
79 composite genotypes due of admixture among animals of different and unclear origin [20]. Here

1 80 we present the first attempt to survey the geo-referenced genomic diversity in wild Chinese RM  
2 81 populations, which is the largest extant population of the species. The current effective  
3 82 population size of Chinese and Indian RM was estimated to be approximate 240,000 and 17,000  
4 83 individuals, respectively, indicating that the Chinese RMs are likely to harbor substantially more  
5 84 genomic diversity compared to their Indian conspecifics [13]. Therefore, this population  
6 85 genomic survey of 81 RMs originating from 17 wild locations across China including  
7 86 phylogenetic and demographic analyses, as well as genome-wide selection scans, corresponds to  
8 87 the most comprehensive characterization of RM genetic diversity to date and aimed at  
9 88 characterizing the processes leading to the extant patterns of variability, as well as identifying the  
10 89 potential implications for the use of these populations in biomedical research.  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22

## 23 91 **Results and Discussion**

### 24 92 **Genetic diversity, phylogeny and population structure**

25 93 Blood and tissue samples from 79 wild-born RMs, representing five subspecies [21, 22], were  
26 94 collected at 17 sites in China (*M. m. tcheliensis*: TH; *M. m. littoralis*: AH, FJ, HB, GX, GZ; *M. m.*  
27 95 *brevicaudus*: HN; *M. m. lasiotis*: SX, SC1, SC2, SC3, SC4; *M. m. mulatta*: YN1, YN2, YN3, YN4,  
28 96 YN5; Fig. 1a). Genome sequences of two additional Chinese RMs (CR1 and CR2) were retrieved  
29 97 from NCBI [9, 23, 24]. Re-sequencing was at a high average depth of  $28.06 \pm 5.08 \times$  for ten  
30 98 individuals and a moderate average depth of  $9.98 \pm 1.05 \times$  for the remainder ( $n=71$ ), with an overall  
31 99 average genome coverage of 93.77% of the RM reference (Mmul\_8.0.1, Supplementary Table 1).  
32 100 A total of 52,534,348 autosome SNPs were identified in these 81 wild Chinese RMs  
33 101 (Supplementary Table 2), and the nucleotide diversity measured by segregating sites (Watterson's  
34 102  $\theta$ ,  $\theta_w$ ) and mean pairwise differences ( $\theta\pi$ ) is 0.00375 and 0.00247, respectively (Table 1). The  
35 103 number of SNPs (all positions with differences to the genome reference) per individual ranged  
36 104 from 7.0 to 9.2 M (mean of 8.50 M; Supplementary Fig. 1 and Supplementary Table 3). Among  
37 105 all detected SNPs, 8,171,139 were shared among all subspecies and 22,768,395 were shared by at  
38 106 least two subspecies, with the remaining SNPs confined to a single subspecies (Supplementary Fig.  
39 107 2a). For each subspecies, the subspecies-specific SNPs (ssSNPs) ranged from 702,099 to  
40 108 7,736,924 and the non-synonymous ssSNPs varied from 3,056 to 25,960 (Supplementary Fig. 2a,  
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

109 b). Among Chinese RM subspecies, *M. m. mulatta* had the highest heterozygosity  
110 ( $2.29 \times 10^{-3} \pm 3.24 \times 10^{-5}$ ), followed by *M. m. lasiotis* ( $2.04 \times 10^{-3} \pm 1.40 \times 10^{-4}$ ) and *M. m. littoralis*  
111 ( $2.00 \times 10^{-3} \pm 1.18 \times 10^{-4}$ ). The lowest heterozygosity rates were found in *M. m. brevicaudus*  
112 ( $1.82 \times 10^{-3} \pm 1.28 \times 10^{-4}$ ) and *M. m. tcheliensis* ( $1.46 \times 10^{-3} \pm 2.65 \times 10^{-4}$ ) (Supplementary Fig. 3).

113 We reconstructed a neighbor-joining (NJ) tree for Chinese RMs based on autosomal SNPs,  
114 using Indian RMs and *M. sylvanus* as outgroups (Fig. 1b and Supplementary Fig. 4). Individuals  
115 from *M. m. lasiotis*, *M. m. brevicaudus* and *M. m. tcheliensis* form monophyletic lineages  
116 respectively, while *M. m. mulatta* and *M. m. littoralis* are paraphyletic. Next, we performed a  
117 population structure analysis using STRUCTURE (version 2.3.4) [25], which estimates  
118 individual ancestry and admixture proportions assuming  $K$  ancestral populations. Plots of  $\Delta K$   
119 generated from STRUCTURE results indicated five genetic clusters present in the full data set  
120 (Fig. 1b and Supplementary Fig. 5). A principal component analysis (PCA) corroborated the  
121 division of Chinese RMs into five groups. The first eigenvector separated *M. m. mulatta* and *M.*  
122 *m. lasiotis* from *M. m. tcheliensis*, *M. m. littoralis* and *M. m. brevicaudus* (variance explained =  
123 7.24%, Tracy-Widom  $P = 4.78 \times 10^{-44}$ ), and the second eigenvector further separated *M. m.*  
124 *tcheliensis*, *M. m. littoralis* and *M. m. brevicaudus* (variance explained = 5.69%, Tracy-Widom  $P$   
125 =  $4.21 \times 10^{-27}$ ) (Fig. 1c, Supplementary Table 4). The division of Chinese RMs into five  
126 geographic lineages supports the former taxonomic division of Chinese RMs into five subspecies  
127 [21, 22]. *M. m. mulatta* (YN1-5) and *M. m. lasiotis* (SC1-4, SX) form the pan-western  
128 populations of Chinese RMs, with both subspecies inhabiting the montane Tibetan Plateau  
129 regions with an altitude  $\geq 1500$  meters above sea level in western China and separated from each  
130 other by the Yangtze River. *M. m. littoralis* (AH, FJ, HB, GX, GZ), *M. m. tcheliensis* (TH) and  
131 *M. m. brevicaudus* (HN) occur in the eastern coastal lowland of China and form the pan-eastern  
132 population. *M. m. tcheliensis* from the Taihang Mountains area (TH) is the northernmost  
133 ( $34^{\circ}54' - 35^{\circ}16' N$ ;  $112^{\circ}02' - 112^{\circ}52' E$ ), while *M. m. brevicaudus*, restricted to Hainan Island, is  
134 the most southern Chinese RM subspecies.

135

### 136 **Demographic and phylogeographic history**

137 The estimated effective population sizes, based on  $\theta_w$  and  $\theta_\pi$  are approximately 93,750 and  
138 61,750 for Chinese RMs (Table 1). In order to infer the ancient demographic history of Chinese

1 139 RMs, we applied a pairwise sequential Markovian coalescent (PSMC) [26] analysis using ten  
2 140 RM individuals with an average sequencing coverage depth higher than 20× (one individual of *M.*  
3  
4 141 *m. tcheliensis* and one of *M. m. brevicaudus*, two of *M. m. lasiotis*, three of *M. m. littoralis* as  
5  
6 142 well as three individuals of *M. m. mulatta*). The inferred PSMC trajectories were very similar for  
7  
8 143 all analyzed individuals throughout most of the species' history reflecting the species'  
9  
10 144 cohesiveness (Fig. 2a). The ancient demographic history of RMs is marked by population  
11  
12 145 fluctuations following the glacial periods during the Pleistocene [27]. Approximately 1,200-800  
13  
14 146 kya all Chinese RMs experienced a population reduction at the time of the Xixiabangma  
15  
16 147 Glaciation (XG), followed by an expansion during the Mid-Pleistocene inter-glaciation (800-200  
17  
18 148 kya). This expansion was then interrupted by the Penultimate Glaciation (PG, 200-130 kya)  
19  
20 149 when suitable habitat might have been lost leading to a population decline [27]. PSMC analyses  
21  
22 150 also suggested that, all the Chinese RMs had a population expansion during the last interglacial  
23  
24 151 (around 100 kya) and a subsequent bottleneck during the Last Glaciation (LG, 70-10 kya) (Fig.  
25  
26 152 2a). Interestingly, the demographic inference by Xue et al. [6] derived from genomic data of a  
27  
28 153 single Chinese RM (CH\_37945) from AH (*M. m. littoralis*) qualitatively resembled the  
29  
30 154 demographic trajectory of *M. m. littoralis* presented herein.

31  
32  
33 155 To further describe the divergence process among the five Chinese RM subspecies, we also  
34  
35 156 employed the SVDquartes approach [28-31] that takes incomplete lineage sorting into account.  
36  
37 157 The obtained phylogenetic tree suggests a “step-by-step” divergence of the five subspecies.  
38  
39 158 Accordingly, the *M. m. mulatta* lineage diverged from that of the remaining Chinese RMs firstly  
40  
41 159 and then the *M. m. lasiotis* diverged from the ancestral lineage of pan-eastern RMs (*M. m.*  
42  
43 160 *tcheliensis*, *M. m. littoralis* and *M. m. brevicaudus*). Subsequently, *M. m. brevicaudus* diverged  
44  
45 161 from the ancestor of *M. m. tcheliensis* and *M. m. littoralis*, the divergence of which occurred lastly  
46  
47 162 (Supplementary Fig. 6). Under this “step-by-step” divergence scenario, we performed the joint site  
48  
49 163 frequency spectrum (SFS) based approach implemented in *fastsimcoal2* [32] to model  
50  
51 164 demographic fluctuations, respective divergence times and gene flow events among the five RM  
52  
53 165 subspecies. Following the divergence between the ancestral lineages of Indian and Chinese RMs  
54  
55 166 (~162 kya), the ancestor of *M. m. mulatta* diverged from the remaining Chinese RMs ~125.8 kya  
56  
57 167 (95 % CI: 92.0-162.1 kya) (Fig. 2b) [6, 13]. Subsequently, *M. m. lasiotis* diverged from the  
58  
59 168 ancestral lineage of pan-eastern RMs ~104.1 kya (95 % CI: 50.2-154.5 kya) near the end of the

169 last interglacial. The divergence time between *M. m. brevicaudus* and the ancestor of *M. m.*  
170 *tcheliensis* and *M. m. littoralis* was estimated at ~61.7 kya (95 % CI: 43.6-115.1 kya), while the  
171 divergence between the latter two occurred ~ 51.3 kya (95 % CI: 7.2-55.4 kya) during the last  
172 glacial maximum [33,34]. Interestingly, the coalescence analysis revealed a large ancestral  
173 population size of the Chinese RMs 125.8 kya (95 % CI: 92.0-162.1 kya) before and a subsequent  
174 population decline and the divergence among the five subspecies (Fig. 2b), which coincided with  
175 the population expansion during the last interglacial (around 100 kya) and the subsequent  
176 bottleneck of Chinese RMs during the Last Glaciation (LG, 70-10 kya) revealed by PSMC  
177 analyses. Our results indicate substantial gene flow occurred between all five extant lineages of  
178 Chinese RMs (Fig. 2b, Supplementary Table 5 and Supplementary Fig. 7).

179 A previous study of mitochondrial DNA identified two major haplogroups dividing Chinese  
180 RMs into a western and an eastern clade. Modern Chinese RMs were thought to have undergone a  
181 northward expansion while entering China via two possible routes: the first into the western  
182 mountains and the second following the eastern coast [35]. Our evolutionary model, however,  
183 suggests a “step-by-step” colonization process of RMs in China (Fig 2c). After the divergence  
184 from the Indian population (~162 kya) [6, 13], the ancestor of Chinese RMs colonized the Tibetan  
185 Plateau from southwestern China, and then experienced a range expansion north and eastwards.  
186 The pan-western population (*M. m. mulatta* and *M. m. lasiotis*) inhabited the western montane  
187 region in China, while the pan-eastern population (*M. m. tcheliensis*, *M. m. littoralis* and *M. m.*  
188 *brevicaudus*) entered the eastern coastal region. These five subspecies further diverged from each  
189 other during the bottleneck caused by the Last Glaciation. Additionally, barriers such as the Yellow,  
190 Yangtze and Pearl rivers and open sea (Fig. 1a) led to further differentiation by limiting gene flow  
191 among them. Water bodies and mountains could therefore be described as driving the formation of  
192 a habitat ‘lattice’ with the different subspecies of RMs occupying different grids in the lattice.

193

#### 194 **Signatures of selection and local adaptation**

195 The wide distribution of Chinese RMs and their respective contrasting habitat types, as well as  
196 their wide use in biomedical studies, makes them an important case study for the analysis of  
197 signatures of local adaptation to divergent selective pressures [36-38]. We identified putative  
198 targets of selection by carrying out pair-wise comparisons between RM subspecies inhabiting the

199 most different environments to increase the chance of finding selection signatures, i.e., *M. m.*  
200 *tcheliensis* that occurs in the northernmost range of the species under cold conditions, and *M. m.*  
201 *brevicaudus* that inhabits the southernmost range of the species, a tropical island. For each  
202 analysis, we compared the five subspecies using the fixation index ( $F_{ST}$ ) and genetic diversity  
203 ( $\theta_\pi$ ), calculated on 50kb long sliding windows (Fig. 3 and Supplementary Figs. 8-13). The top 5%  
204 of the windows with the largest  $F_{ST}$  and  $\theta_\pi$  ratios ( $\theta_{\pi 2} / \theta_{\pi 1}$ ) in each pair-wise comparison were  
205 considered to be potentially under positive selection. For each subspecies, we identified the  
206 intersection of potential selective-sweep regions generated by all the pair-wise comparisons  
207 between a subspecies and each of the other subspecies (four pairwise comparisons in each case)  
208 (Supplementary Fig. 8). We used these consistent selective-sweep regions for further analyses, as  
209 they represent robust putative positively selected regions. The sizes of candidate selective-sweep  
210 regions ranged from 0.100 Mb to 11.075 Mb and the number of genes located in these regions,  
211 which are expected to represent targets of selection for each subspecies, varied from 6 to 176 in  
212 different subspecies (Supplementary Table 6).

213 *M. m. tcheliensis* from the Taihang (TH) Mountains area is the northernmost population of  
214 the species. The TH Mountains are characterized by a continental monsoon climate, and  
215 conditions for RMs are harsh during winter and early spring with extreme cold temperatures of –  
216 14°C [39]. Food resources are limited and consist mainly of barks, twigs, roots of crops and  
217 withered grass, thus, all sources are high in fiber, but low in energy and nutritional value [40, 41].  
218 Therefore, *M. m. tcheliensis* suffers from starvation due to food shortage during winter and early  
219 spring. In starvation, blood glucose levels are maintained by gluconeogenesis through which  
220 glucose are converted from other molecules, such as amino acids and lactic acid [42]. For *M. m.*  
221 *tcheliensis*, the positive selection genes are enriched in the gene ontology (GO) term “fructose 1,  
222 6-bisphosphate 1-phosphatase activity” with two genes (*Fbp1*, *Fbp2*, modified Fisher Exact  
223  $P=1.90E-02$ ; Fig. 3c, d; Supplementary Table 7). These two genes encode for fructose-1,  
224 6-bisphosphatase 1 and fructose-1, 6-bisphosphatase isozyme 2 which catalyze the hydrolysis of  
225 fructose 1, 6-bisphosphate and play a rate-limiting role in gluconeogenesis. Furthermore, in  
226 starved zebrafish it was shown that the expression of *Fbp1* was significantly upregulated in brain  
227 and liver tissues [43]. The positive selection genes are also enriched in other terms and pathway  
228 related to gluconeogenesis, including KEGG pathway “Fructose and mannose metabolism”

1 229 (modified Fisher Exact  $P=4.35E-02$ ) and GO terms “hexose biosynthetic process”,  
2 230 “monosaccharide biosynthetic process” and “cellular carbohydrate biosynthetic process”  
3  
4 231 (modified Fisher Exact  $P=3.36E-02$ ,  $P=4.64E-02$  and  $P=2.65E-02$ ; Supplementary Table 7). Our  
5  
6 232 findings suggest that the regulation of gluconeogenesis might be a mechanism of *M. m. tcheliensis*  
7  
8 233 to adapt to food shortage in winter.

9  
10 234 According to Bergman’s rule, animals living in cold climates tend to have larger body sizes  
11  
12 235 compared to their relatives in warm climates (i.e. they have a lower surface area to volume ratio),  
13  
14 236 so they radiate less body heat per unit of mass [44]. Consistent with this expectation, among all  
15  
16 237 RM subspecies, *M. m. tcheliensis* exhibits the largest body size and mass, and the largest head and  
17  
18 238 chest circumference (Fig. 3b and Supplementary Table 8) [40, 45]. Among the consistent  
19  
20 239 signatures of positive selection identified in *M. m. tcheliensis* (176 genes), we found signatures of  
21  
22 240 selective sweeps in eight genes linked to limb morphogenesis or skeletal system development  
23  
24 241 (Supplementary Table 6). Among these genes, *Fto* and *Rpgrip1l* play an essential role in postnatal  
25  
26 242 growth of mammals [46]. Mice lacking *Fto* completely display immediate postnatal growth  
27  
28 243 retardation with shorter body length, lower body weight, and lower bone mineral density than  
29  
30 244 control animals [47]. Furthermore, *Sox5* and *Sox6* (Fig. 3c, d) play an essential role in synovial  
31  
32 245 joint morphogenesis via promoting both growth plate and articular chondrocyte differentiation  
33  
34 246 [48]. Mutations in *Atp6v0a4* could cause developmental delay and delayed closure of the anterior  
35  
36 247 fontanelle in human [49], while expression of *Ext2* enhances the bone formation in mice [50]  
37  
38 248 These genes involved in the growth and development of the skeletal system and appendages are  
39  
40 249 likely contributors to the larger body size of *M. m. tcheliensis*, and represent an undescribed  
41  
42 250 adaptive pathway for primates living in colder climates.

43  
44  
45 251 In contrast, *M. m. breviceaudus* inhabits the tropical island of Hainan (HN) where it copes  
46  
47 252 with a mean annual temperature of 24°C. *M. m. breviceaudus* has the smallest body size, the  
48  
49 253 smallest body mass, and the shortest tail among RM subspecies [45]. As described above, they  
50  
51 254 radiate more body heat per unit of mass (Bergman’s rule) [44]. We found 127 putatively selected  
52  
53 255 genes in *M. m. breviceaudus* (Supplementary Table 6), four of which were found to be enriched in  
54  
55 256 GO term “Bone morphogenetic protein (BMP) signaling pathway” (modified Fisher Exact  
56  
57 257  $P=4.65E-02$ ; Supplementary Table 9) and two genes were found to be enriched in GO term  
58  
59 258 “I-SMAD binding (modified Fisher Exact  $P=4.65E-02$ ; Supplementary Table 9)”. BMP and



1 259 I-SMAD signaling pathways are involved in the development of bones and the skeleton [51, 52].  
2 260 Mutations in *Axin1*, a gene of the I-SMAD pathway, cause kinked tails in mice [53]. In *M. m.*  
3 *brevicaudus*, we found two non-synonymous mutations in this gene (A674G, T656I)  
4 261  
5  
6 262 (Supplementary Fig. 14 and Supplementary Table 10, 11).

7  
8 263 Additionally, putatively selected genes in *M. m. brevicaudus* (Fig. 3c, d, Supplementary  
9  
10 264 Table 6) were also involved in GO terms related to cardiovascular system and blood circulation.  
11  
12 265 For example, *Aggf1* related to GO term “blood vessel morphogenesis” and *Ctnna3* related to GO  
13  
14 266 term “regulation of heart rate by cardiac conduction”. The up-regulated *Aggf1* expression is  
15  
16 267 capable of increasing blood flow in mouse hindlimb [54]. In addition, *Hspa4*, heat shock 70kDa  
17  
18 268 protein 4, is directly involved in GO term “response to temperature stimulus”. We thus  
19  
20 269 hypothesize that the cardiovascular system of *M. m. brevicaudus* might play an important role in  
21  
22 270 stabilizing body temperature, assisted by blood flow through different body parts requiring good  
23  
24 271 fluidity and vascular permeability to transfer heat out of the body [55]. Testing these hypotheses  
25  
26 272 needs further functional assays, however, these genes, together with the positively selected genes  
27  
28 273 identified in *M. m. tcheliensis*, are known to be relevant to human physical function, and thus are  
29  
30 274 likely of importance in the adaptation of Chinese RMs to different climate conditions.

31  
32 275 Both coding and non-coding changes could contribute to local adaptations of organisms [56].  
33  
34 276 To further investigate the adaptive mechanism of *M. m. tcheliensis* and *M. m. brevicaudus* to the  
35  
36 277 opposite climates (cold versus hot), we focused on SNPs in the gene regions of above described  
37  
38 278 candidate genes. A total of 5817 SNPs were found with significant differences at the 5% level in  
39  
40 279 the distributions of genotypes between these two subspecies, and 10 SNPs were non-synonymous  
41  
42 280 variants (Supplementary table 10 and 11). In *M. m. tcheliensis*, non-synonymous mutations were  
43  
44 281 found in the coding regions of *Atp6v0a4* (R667Q), *Ext2* (I363M), *Fto* (N10S) and *Rpgrip11*  
45  
46 282 (R1281Q) (Supplementary table 11 and Supplementary Fig. 14), implying that selection might  
47  
48 283 have acted on protein sequence changes. No non-synonymous changes were detected in *Fbp1*,  
49  
50 284 *Fbp2*, *Sox5* and *Sox6*. However, SNPs are located in the 1kb up/downstream, 5' and 3' UTR, and  
51  
52 285 intronic regions of these genes (Supplementary table 10), indicating selection on non-coding  
53  
54 286 regulatory variants. Correspondingly, non-synonymous mutations in *Aggf1* (H343Y), *Axin1*  
55  
56 287 (A674G, T656I), *Hspa4* (I782V) and *Ctnna3* (V551I, T577M) were revealed for *M. m.*  
57  
58 288 *brevicaudus* (Supplementary table 11 and Supplementary Fig. 14).

1 289 Besides the genes related to the adaptation to various climate conditions, we also found  
2 290 signatures of positive selection in genes related to the nervous system. In *M. m. tcheliensis* the  
3  
4 291 176 identified candidate genes are enriched in GO term “synapse” (modified Fisher Exact  
5  
6 292  $P=4.28E-02$ ; Supplementary Table 7) with eight genes, and two of these gene, *Gabra2* and  
7  
8 293 *Chrm2* are associated with alcohol dependence [57]. For *M. m. brevicaudus*, 18 putatively  
9  
10 294 selected genes related to nervous system development were found. For example, *Dcc* is reported  
11  
12 295 to be required for long-term potentiation and memory [58]. *Auts2*, one of the eight putatively  
13  
14 296 selected genes in *M. m. lasiotis*, has been shown to regulate neuronal migration, and mutations in  
15  
16 297 this gene cause mental dysfunction in human [59] (Supplementary Table 6). Our findings suggest  
17  
18 298 that RM subspecies have experienced different adaptative processes in the nervous system and  
19  
20 299 respective genomic differences should be taken into account when animals are selected for  
21  
22 300 neurobiological research.  
23  
24

25 301

### 27 302 **Disease-causing variants and implication for biomedical research**

28  
29 303 Given the large evolutionary similarity between macaques and humans, human diseases are  
30  
31 304 better modeled in RMs than in many other animals. Thus, variants in RMs that match to  
32  
33 305 orthologous human variants annotated as ‘pathogenic’ are of particular interest. We examined  
34  
35 306 presumed homologous Chinese RM SNPs in the human genome and a total of 34,850,330 RM  
36  
37 307 SNPs analyzed in this study were successfully identified in the human genome (hg19). Among  
38  
39 308 these SNPs, 118 variants matched human variants with the accordant reference alleles and  
40  
41 309 alternative alleles were annotated as ‘disease causing’ in HGMD or pathogenic in ClinVar. These  
42  
43 310 118 RM SNPs affect genes that cause specific human diseases including acromesomelic  
44  
45 311 dysplasia maroteaux type, anonychia, atranferrinemia, blau syndrome, Carcinoma of colon,  
46  
47 312 Charcot-Marie-Tooth disease, deafness, early infantile epileptic encephalopathy 7, glycogen  
48  
49 313 storage disease and others (Supplementary Table 12). Among these 118 SNPs, only seven  
50  
51 314 pathogenic SNPs are shared by all five subspecies, while 82 are subspecies-specific (Fig. 4c,  
52  
53 315 Supplementary Table 12). For example, the SNP rs116229331 in the gene *Unc13d* (human Chr17:  
54  
55 316 73836585C>T), known to cause juvenile idiopathic arthritis in humans [60], has a RM  
56  
57 317 homologue (RM Chr16: 69559126 C>T, Fig. 4a) that is present in *M. m. tcheliensis*, *M. m.*  
58  
59 318 *brevicaudus* and *M. m. littoralis*, but absent in *M. m. lasiotis* and *M. m. mulatta*. Another  
60  
61

1 319 pathogenic variant (rs397514345, human Chr3: 15686724 A>C) in the *Btd* gene is involved in  
2 320 biotinidase deficiency [61]. Its homologous RM variant (RM Chr2: 172277927 A>C, Fig. 4a) is  
3  
4 321 found only in *M. m. lasiotis* and *M. m. mulatta*. In addition, we also identified 16  
5  
6 322 non-synonymous SNPs in the *Noca3* gene, which encodes a protein that modulates the  
7  
8 323 replication and transcriptional reactivation of HIV-1 during virus latency [62] (Fig. 4b). Ten of  
9  
10 324 these 16 non-synonymous SNPs are private to one subspecies (Supplementary Table 13). The  
11  
12 325 effects of these variants on HIV-1 replication and reactivation are unknown and need further  
13  
14 326 investigation, but the high number of mutations suggests a complex response of the host to the  
15  
16 327 virus.

18  
19 328 Overall, these findings suggest that the genomic architecture of Chinese RMs used in  
20  
21 329 biomedical research and their geographic origin could strongly influence the outcome of  
22  
23 330 biomedical experiments and should be taken into account when using Chinese RMs in clinical and  
24  
25 331 neurobiological research. Unfortunately, genome wide screening of RMs used in biomedical  
26  
27 332 research is so far only rarely conducted and uncharacterized animals are most often used.  
28  
29 333 Importantly, individuals from all five Chinese RM subspecies are used in biomedical research [63,  
30  
31 334 64]. Combined with our data, nine of the 26 captive Chinese RMs reported by Zhong *et al.* [7]  
32  
33 335 were found to cluster with *M. m. littoralis*, 16 with *M. m. lasiotis* and one with *M. m. mulatta* (Fig.  
34  
35 336 4d). Thus, the data and results presented here provide the basis to trace the origin of captive RMs  
36  
37 337 and to allow for the selection of appropriate animal models when testing for particular diseases,  
38  
39 338 and are thus a significant contribution to the “3Rs” principle, which aims to reduce, refine, and  
40  
41 339 replace experimental animals [65].  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

## 340 **Conclusion**

341 We present the first description of the evolutionary history and genomic variation of  
342 geo-referenced wild RMs throughout China, including scenarios on potential functions of this  
343 variation in adaptation to local environments. This genomic resource represents a valuable  
344 contribution to the understanding of the biology and evolution of a highly successful and  
345 important biomedical research species. In particular, it is important to note that due to the  
346 difference in evolutionary history of the subspecies identified here, it can be expected that  
347 animals originating from different regions may react differently to experimental tests, and thus  
348 their background needs to be assessed beforehand [10]. Our results highlight the importance that  
349 genome typing can play in biomedical research where animal origins are uncertain, and the  
350 resources generated here provide a baseline for genomic assessment of biomedical research  
351 populations, genetic resource conservation and for refined usage of RMs in future research.

## 352 **Materials and Methods**

### 353 **Ethics statement**

354 The methods were carried out in accordance with the approved guidelines of the Good  
355 Experimental Practices adopted by the Institute of Zoology, Chinese Academy of Sciences  
356 (CAS). All experimental procedures and animal collection were conducted under the supervision  
357 of the Committee for Animal Experiments of the Institute of Zoology, Chinese Academy of  
358 Sciences.

### 359 **Sample Collection and Sequencing**

360 Samples from 79 individuals with information about geographic origin were collected from 17  
361 local wildlife rescue center, which covered most of the species' range in China. Muscle samples  
362 were collected from deceased individuals and the blood samples were taken during routine  
363 physical examinations. Total genomic DNA was extracted from blood or tissue samples using  
364 standard phenol/chloroform methods. For each individual, ~3 µg DNA was sheared into  
365 fragments of 500 bp with the Covaris system. DNA fragments were then processed and  
366 sequenced using the Illumina HiSeq 2000 and 2500 platform. Furthermore, published genomic  
367 data for two individuals were download form NCBI [9,23] and filtered using the same conditions.  
368 Raw reads were first filtered with the following criteria: (1) reads with unidentified nucleotides  
369 (N) exceeded 10% were discarded, (2) reads with the proportion of low quality base (phred  
370 quality <=5) larger than 50% were discarded. After the quality control, a total of 3,095.6 Gb of  
371 high quality sequences with 22.53 billion pair-end reads (100 or 125 bp) were generated.

### 373 **Sequence Data Pre-processing and Variant Calling**

374 High-quality sequence reads were mapped to the macaque reference genome, Mmul\_8.0.1 [66],  
375 using the Burrows–Wheeler Aligner 0.7.10-r789 (BWA, RRID:SCR\_010910) [67]. Sequence  
376 Alignment/Map (SAM) format files were imported to SAMtools v0.1.19 (SAMtools,  
377 RRID:SCR\_002105) [68] for sorting and then imported to Picard version 1.118 (Picard,  
378 RRID:SCR\_006525) (<http://broadinstitute.github.io/picard/>) for removing duplicated reads. To  
379 improve the quality of sites reported, we performed SNP calling following GATK's best practice,  
380 version 3.3–0 (GATK, RRID: SCR\_001876) on autosomal sites only [69]. We get the GVCF file

1 381 for each individual using the “HaplotypeCaller” method in GATK and then using  
2 382 GenotypeGVCFs-based method with the “-includeNonVariantSites” flag to get the population  
3 383 VCF file including all the confident sites. After that, we first applied the “SelectVariants” to  
4 384 exclude the Indel and split the variant and non-variant sites. Then we applied the hard filter  
5 385 command ‘VariantFiltration’ to exclude potential false-positive variant calls with the following  
6 386 criteria: “-filterExpression ‘QD < 5.0 | FS > 60.0 | MQ < 40.0 | ReadPosRankSum< -8.0 ||  
7 387 MQRankSum < -12.5’” and “--genotypeFilterExpression ‘DP < 4.0’”. Additionally, the sites are  
8 388 filtered out if there is a 'N' is in the reference sequence; if the site is fixed difference to the  
9 389 reference genome or if the site including more than 20% missing genotypes. For non-variant  
10 390 sites, we filtered the sites if there is a 'N' is in the reference sequence or if the site including more  
11 391 than 20% missing genotypes. All the SNPs were annotated by ANNOVAR v2013-06-21  
12 392 (ANNOVAR, RRID:SCR\_012821) [70] (Supplementary Table 2). For each individuals, the  
13 393 heterozygosity was calculated as heterozygous SNP rate across the whole genome based on the  
14 394 whole number of sites that are callable (Supplementary Table 3).

### 31 396 **Genetic Diversity and Structure Analysis**

32 397 A neighbor-joining (NJ) tree was constructed for the 81 individuals based on the autosomal  
33 398 genome data using the software TreeBeST. The bootstrap was set to 1,000 times to assess branch  
34 399 support, with the genome information of Indian RMs and *M. sylvanus* as outgroups. FigTree  
35 400 (<http://tree.bio.ed.ac.uk/software/figtree/>, v1.4.0) was used to visualize the phylogenetic tree (Fig.  
36 401 1b and Supplementary Fig. 4). Population structure analysis was performed using the software  
37 402 STRUCTURE 2.3.4 [25], which estimates individual ancestry and admixture proportions  
38 403 assuming  $K$  ancestral populations. We ran STRUCTURE five times to assess convergence and  
39 404 tested the number of genetic clusters ( $K$ ) from 2-9 (Supplementary Fig. 5). We also carried out a  
40 405 principle component analysis (PCA) using the smartPCA program from the Eigensoft package,  
41 406 v5.0 (Eigensoft, RRID:SCR\_004965) [71]. To determine the significance level of principal  
42 407 components, a Tracy-Widom test was done after the PCA (Supplementary Table 4). Decay of  
43 408 linkage disequilibrium against physical distance for the different populations was calculated  
44 409 using the Haploview software [72] with the maxdistance set as 500kb (Supplementary Fig. 15).

60 410

## 411 Demographic and Divergence Inference Using PSMC and Fastsimcoal2

412 We called the consensus sequences using Samtools mpileup [68] by applying: “samtools mpileup  
413 -q 1 -C 50 -S -D -m 2 -F 0.002 -u -f \*.fa(genome) \*.bam | bcftools view -c - | vcfutils.pl vcf2fq -d  
414 10 -D 100 -Q 20 - > \*.psmc.fq” and “fq2psmcfa -q10 -s 100 \*.psmc.fq >\*.psmc.fa” . To ensure  
415 the quality of consensus sequences, we used data of ten individuals with an average  
416 coverage >20× (22.20-34.32×). The PSMC model [26] was used to estimate the population  
417 histories from the individual genomes (sex chromosomes excluded) with the following parameters:  
418 -N30 -t15 -r5 -p ‘4+25×2+4+6’. We chose a generation length of 11 years and a mutation rate  
419 per generation ( $\mu$ ) of  $1.0 \times 10^{-8}$  (for the rationale to use these values see [6, 73]).

420 We used PAUP\* 4.0a142 (PAUP, RRID:SCR\_014931) [30] to run SVDquartets to estimate  
421 the branching pattern among the five subspecies with the following command: SVDQuartets  
422 SpeciesTree=yes bootstrap evalQuartets=all seed=0 nthreads=40. The joint site frequency  
423 spectrum (SFS) approach implemented in *fastsimcoal2* [32] was performed to model more recent  
424 demographic fluctuations and respective divergence times based on the species tree estimation by  
425 SVDquartets. VCF file containing callable variant sites was used converted to fastsimcoal style  
426 folded SFS. To mitigate the effect of linkage disequilibrium, we filtered out the SNPs located  
427 within 10 kb from genes and then we took one SNPs every 10kb randomly. The multidimensional  
428 folded SFS for all the five subspecies is generated by easySFS  
429 (<https://github.com/isaacovercast/easySFS#easysfs>). During likelihood calculation, a conditional  
430 maximization algorithm (ECM) is used to maximize the likelihood of each parameter while  
431 keeping the others stabilized. This ECM procedure runs through 40 cycles where each  
432 composite-likelihood was calculated using 100,000 coalescent simulations. Additionally, in order  
433 to avoid likelihood estimates that oversample parameter values at local maxima across the  
434 composite likelihood surface, we ran 50 replicates with each starting from different initial  
435 conditions. We chose the replicate with the highest estimated maximum likelihood score to  
436 estimate confidence intervals using parametric bootstrapping. The SFS used in bootstrap was  
437 simulated with the parameter values from the highest likelihood model and then new parameter  
438 values re-estimated from the simulated SFS. We ran 100 parametric bootstraps (Supplementary  
439 Fig. 7).

440

## 441 **Positive Selection**

442 To identify genomic regions that may have been subject to selection for each subspecies  
443 inhabited in different habitats, we scanned the genome using one-to-one pair-wise comparisons  
444 between all five subspecies. We calculated the genome-wide distribution of  $F_{ST}$  values [74] and  
445  $\theta_\pi$  ratios for each pairwise comparison among five RM subspecies. We calculated  $\theta_\pi$  for each  
446 population and the  $F_{ST}$  between the two populations in each comparison using VCFtools  
447 (VCFtools, RRID:SCR\_001235) [75] with a genome-wide sliding window strategy (50-kb in  
448 length with 25-kb step). The  $F_{ST}$  values were Z-transformed and the log value of  $\theta_\pi$  ratio ( $\theta_{\pi 2}$   
449 / $\theta_{\pi 1}$ ) was estimated. Candidate regions under positive selection were extracted based on the top 5%  
450 of log-odds ratios for both Z ( $F_{ST}$ ) and log ( $\theta_\pi$ -ratio). Finally, for each subspecies we used the  
451 intersection of putatively selected regions generated by all the pair-wise comparisons with other  
452 subspecies as the candidate regions under positive selection (i.e. consistent signatures of  
453 selective sweeps). Genes located in these regions are expected to represent targets of selection.  
454 Functional classification and enrichment analysis of GO categories and KEGG pathways for  
455 these candidate genes were performed using DAVID v6.8 (DAVID, RRID:SCR\_001881) [76].  
456 The modified Fisher Exact  $P$ -value cut off was 0.05. Chi-square and  $P$ -values for the allele  
457 frequencies in *M. m. tcheliensis* vs. *M. m. brevicaudus* for the re-sequenced SNPs from the  
458 candidate genes were assessed with the Haploview program [72].

## 460 **Genomic divergence and implication for biomedical research**

461 A total of 118 out of 52,534,348 RM SNPs analyzed in this study were successfully mapped to  
462 human reference sequence version hg19 (GRCh37) using liftOver  
463 (<https://genome.ucsc.edu/cgi-bin/hgLiftOver>) and were annotated as ‘disease causing’ in HGMD  
464 (version 2015.1) or pathogenic in ClinVar (downloaded 25/02/2018) (Supplementary Table 12).

## 466 **Data Availability**

467 All data generated from this study have been submitted to the NCBI Sequence Read Archive  
468 (SRA) under BioProject PRJNA345528. The datasets supporting the results of this article are  
469 available in the *GigaScience* GigaDB repository [77].



1 470 **Competing interests**

2  
3 471 The authors declare that they have no competing interests.  
4  
5  
6 472

7  
8 473 **Acknowledgments**  
9

10 474 This project was sponsored by the following grants: Ming Li (Key Project of National Natural  
11  
12 475 Science Foundation of China, 31530068; Strategic Priority Research Program of the Chinese  
13  
14 476 Academy of Sciences, XDA19050202 and XDB31000000; National Key R&D Program of  
15  
16 477 China, 2016YFC0503200; Creative Research Group Project of NSFC; and Science &  
17  
18 478 Technology Department of Sichuan Province, 2018JZ0008); Zhijin Liu (Natural Science  
19  
20 479 Foundation of China, 31471989). The authors thank Baoguo Li, Meng Yao, Songtao Guo, Jiqi  
21  
22 480 Lu, Zhenlong Wang, Xuelong Jiang, Tao Meng and Qihai Zhou for their help in sampling;  
23  
24 481 Daniel Pitt, Quan Kang, Qi Wu, Chuanyun Li and Qi Pan for their assistance in data analysis.  
25  
26  
27 482

28  
29 483 **Author contributions**  
30

31 484 M. L., Z. L. and M. B conceived the study and designed the project. Z. L., X. T., P. O., X. Z., L.  
32  
33 485 Z. and S. T. managed the project, performed the analyses and wrote the manuscript. Z. L., B. S.  
34  
35 486 and H. X. prepared samples. Z. L., X. T. and P. O. performed genetic analyses. Z. L., X. T., P. O.,  
36  
37 487 B. R., L. Z., G. L., Z. Y., Z. P., Z. X., C. R., M. B. and M. L. discussed the data. Z. L. and X. T.  
38  
39 488 wrote the manuscript with contributions from P. O., B. W., H. X., W. Z., C. R., M. B. and M. L.;  
40  
41 489 all authors contributed to data interpretation.  
42  
43

44 490 **Supplementary Material**  
45

46  
47 491 Supplementary information, figures S1-S15 and tables S1-S13 are available on line.  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

## References

1. Moreno-Estrada A, Gignoux CR, Fernández-López JC et al. Human genetics. The genetics of Mexico recapitulates Native American substructure and affects biomedical traits. *Science* 2013; **344**:1280–1285.
2. Allentoft ME, Sikora M, Sjögren KG et al. Population genomics of Bronze Age Eurasia. *Nature* 2015; **522**:167–172.
3. Sudmant PH, Rausch T, Gardner EJ et al. An integrated map of structural variation in 2,504 human genomes. *Nature* 2015; **526**:75–81.
4. Maestripieri D. *Macachiavellian intelligence: How rhesus macaques and humans have conquered the world*. 2007. The University of Chicago Press, Chicago.
5. Zinner D, Fickenscher GH, Roos C. Family Cercopithecidae (Old World Monkeys). *Handbook of the Mammals of the World*. 2013; Pp. 550-753 in: Mittermeier RA, Rylands AB, Wilson DE. eds. Vol. 3. Primates. Lynx Edicions, Barcelona.
6. Xue C, Raveendran M, Harris RA et al. The population genomics of rhesus macaques (*Macaca mulatta*) based on whole genome sequences. *Genome Res* 2016; **26**:1651–1662.
7. Zhong X, Peng J, Shen QS et al. RhesusBase PopGateway: Genome-Wide Population Genetics Atlas in Rhesus Macaque. *Mol Biol Evol* 2016; **33**:1370–1375.
8. Fawcett GL, Raveendran M, Deiros DR et al. Characterization of single-nucleotide variation in Indian-origin rhesus macaques (*Macaca mulatta*). *BMC Genomics* 2011; **12**:311.
9. Yan G, Zhang G, Fang X et al. Genome sequencing and comparison of two nonhuman primate animal models, the cynomolgus and Chinese rhesus macaques. *Nature Biot* 2011; **29**:1019-1023.
10. Haus T, Ferguson B, Rogers J et al. Genome typing of nonhuman primate models: implications for biomedical research. *Trends Genet* 2014; **30**:482–487.
11. Flynn S, Satkoski J, Lerche N et al. Genetic variation at the TNF-alpha promotor and malaria susceptibility in rhesus (*Macaca mulatta*) and long-tailed (*Macaca fascicularis*) macaques. *Infect Genet Evol* 2009; **9**:769–777.
12. de Groot NG, Heijmans CMC, Koopman G et al. TRIM5 allelic polymorphism in macaque species/populations of different geographic origins: its impact on SIV vaccine studies. *Tissue*

- 521 Antigens. 2011; **78**:256–62.
- 522 13. Hernandez RD, Hubisz MJ, Wheeler DA et al. Demographic histories and patterns of linkage  
523 disequilibrium in Chinese and Indian rhesus macaques. *Science* 2007; **316**:240–243.
- 524 14. Champoux M, Higley JD, Suomi SJ. Behavioral and physiological characteristics of Indian  
525 and Chinese-Indian hybrid rhesus macaque infants. *Dev Psychobiol* 1997; **31**:49–63.
- 526 15. Trichel AM, Rajakumar PA, Murphey-Corb M. Species-specific variation in SIV disease  
527 progression between Chinese and Indian subspecies of rhesus macaque. *J Med Primatol* 2002;  
528 **31**:171–178.
- 529 16. Tosi AJ, Morales JC, Melnick DJ. Paternal, maternal, and biparental molecular markers  
530 provide unique windows onto the evolutionary history of macaque monkeys. *Evolution* 2003;  
531 **57**:1419–1435.
- 532 17. Smith DG. Genetic characterization of Indian-origin and Chinese-origin rhesus macaques  
533 (*Macaca mulatta*). *Comp Med* 2005; **55**:227–230.
- 534 18. Ferguson B, Street SL, Wright H et al. Single nucleotide polymorphisms (SNPs) distinguish  
535 Indian-origin and Chinese-origin rhesus macaques (*Macaca mulatta*). *BMC Genomics* 2007;  
536 **8**:43.
- 537 19. Kubisch HM, Falkenstein KP, Deroche CB et al. Reproductive efficiency of captive Chinese-  
538 and Indian-origin rhesus macaque (*Macaca mulatta*) females. *Am J Primatol* 2012; **74**:174–  
539 184.
- 540 20. Kanthaswamy S, Johnson Z, Trask JS et al. Development and validation of a SNP-based assay  
541 for inferring the genetic ancestry of rhesus macaques (*Macaca mulatta*). *Am J Primatol* 2014;  
542 **76**:1105–1113.
- 543 21. Fooden J. Systematic review of the rhesus macaque, *Macaca mulatta* (Zimmermann, 1780).  
544 *Field Zool* 2000; **96**:1–180.
- 545 22. Jiang X, Wang Y, Ma S. Taxonomic revision and distribution of subspecies of rhesus monkey  
546 (*Macaca mulatta*) in China. *Zool Res* 1991; **12**:241–247.
- 547 23. Fang X, Zhang Y, Zhang R et al. Genome sequence and global sequence variation map with  
548 5.5 million SNPs in Chinese rhesus macaque. *Genome Biol* 2011; **12**:R63.
- 549 24. Prado-Martinez J, Sudmant PH, Kidd JM et al. Great ape genetic diversity and population  
550 history. *Nature* 2013; **499**:471–475.

- 1 551 25. Evanno G, Regnaut S, Goudet J. Detecting the number of clusters of individuals using the  
2 552 software STRUCTURE: a simulation study. *Mol Ecol* 2005; **14**:2611–2620.
- 3  
4 553 26. Li H, Durbin R. Inference of human population history from individual whole-genome  
5 554 sequences. *Nature* 2011; **475**:493–496.
- 6  
7 555 27. Zheng B, Xu Q, Shen Y. The relationship between climate change and Quaternary glacial  
8 556 cycles on the Qinghai–Tibetan Plateau: review and speculation. *Quatern Int* 2002; **97**:93–  
9 557 101.
- 10  
11 558 28. Chifman J, Kubatko L. Identifiability of the unrooted species tree topology under the  
12 559 coalescent model with time-reversible substitution processes, site-specific rate variation, and  
13 560 invariable sites. *J Theor Biol* 2014; **374**:35–47.
- 14  
15 561 29. Chifman J, Kubatko L. Quartet inference from SNP data under the coalescent model.  
16 562 *Bioinformatics* 2014; **30**:3317–3324.
- 17  
18 563 30. Swofford, D, et al; PAUP\*. *Phylogenetic Analysis Using Parsimony (\*and other methods)*.  
19 564 Version 4. Sinauer Associates, Sunderland, Massachussets. 2003.
- 20  
21 565 31. Song S, Liu L, Edwards SV, Wu S. Resolving conflict in eutherian mammal phylogeny using  
22 566 phylogenomics and the multispecies coalescent model. *Proc Natl Acad Sci USA* 2012; **109**:  
23 567 14942–14947.
- 24  
25 568 32. Excoffier, L. Dupanloup I, Huerta-Sánchez E et al Robust demographic inference from  
26 569 genomic and SNP data. *PLoS Genet* 2013; **9**:e1003905.
- 27  
28 570 33. Owen LA, Finkel RC, Caffee MW. A note on the extent of glaciation throughout the Himalaya  
29 571 during the global Last Glacial Maximum. *Quaternary Sci Rev* 2002; **21**:147–157.
- 30  
31 572 34. Owen LA. Latest Pleistocene and Holocene glacier fluctuations in the Himalaya and Tibet.  
32 573 *Quaternary Sci Rev* 2009; **28**:2150–2164.
- 33  
34 574 35. Wu S, Luo J, Li Q et al. Ecological genetics of Chinese rhesus macaque in response to  
35 575 mountain building: all things are not equal. *PLoS ONE* 2013; **8**:e55315.
- 36  
37 576 36. Yi X, Liang Y, Huerta-Sanchez E et al. Sequencing of 50 human exomes reveals adaptation to  
38 577 high altitude. *Science* 2010; **329**:75–78.
- 39  
40 578 37. Bhatia G, Patterson N, Pasaniuc B et al. Genome-wide comparison of African-ancestry  
41 579 populations from CARE and other cohorts reveals signals of natural selection. *Am J Hum*  
42 580 *Genet* 2011; **89**:368–381.

- 1 581 38. Zhao SC, Zheng PP, Dong SS et al. Whole-genome sequencing of giant pandas provides  
2 582 insights into demographic history and local adaptation. *Nat Genet* 2013; **45**:67–71.
- 3  
4 583 39. Tian JD, Wang ZL, Lu JQ, Wang BS, Chen JR. Reproductive Parameters of Female *Macaca*  
5 584 *mulatta tcheliensis* in the Temperate Forest of Mount Taihangshan, Jiyuan, China. *Am J*  
6 585 *Primatol* 2013; **75**:605–612.
- 7  
8 586 40. Zhao X, Zhang H, Lv X et al. Survey and research of morphological characters of monkeys  
9 587 (*Macaca mulatta*) in the Taihang Mountains. *J Henan Nor Uni* 1989; **62**:120–125.
- 10 588 41. Lu JQ, Hou JH, Wang HF, Qu WY. Current status of *Macaca mulatta* in Taihangshan  
11 589 Mountains Area, Jiyuan, Henan, China. *Int J Primatol* 2007; **28**:1085–1091.
- 12  
13 590 42. Sadava DE, Heller HC, Orians GH, Purves WK, Hillis DM. *Life: The Science of Biology*,  
14 591 8th edn. Macmillan, New York; 2008.
- 15  
16 592 43. Drew RE, Rodnick KJ, Settles M et al. Effect of starvation on transcriptomes of brain and  
17 593 liver in adult female zebrafish (*Danio rerio*). *Physiol Genomics* 2008; **35**:283–295.
- 18  
19 594 44. Bergmann C. Über die Verhältnisse der Wärmeökonomie der Thiere zu ihrer Grösse.  
20 595 Göttinger Studien 1847; **3**:595–708.
- 21  
22 596 45. Zhang P, Lyu MY, Wu CF et al. Variation in body mass and morphological characters in  
23 597 *Macaca mulatta brevicaudus* from Hainan, China. *Am J Primatol* 2016; **78**:679–698.
- 24  
25 598 46. Jevsinek Skok D, Kunej T, Kovac M et al. FTO gene variants are associated with growth  
26 599 and carcass traits in cattle. *Animal Genetics* 2016; **47**: 219–222.
- 27  
28 600 47. Gao X, Shin YH, Li M, Wang F, Tong Q, Zhang PM. The Fat Mass and Obesity Associated  
29 601 Gene FTO Functions in the Brain to Regulate Postnatal Growth in Mice. *PLoS ONE* 2010; **5**:  
30 602 e14005.
- 31  
32 603 48. Dy P, Smits P, Silvester A et al. Synovial joint morphogenesis requires the chondrogenic  
33 604 action of Sox5 and Sox6 in growth plate and articular cartilage. *Dev Biol* 2010; **341**:346–  
34 605 359.
- 35  
36 606 49. Greally MT, Kalis NN, Agab W et al. Autosomal recessive cutis laxa type 2A (ARCL2A)  
37 607 mimicking Ehlers - Danlos syndrome by its dermatological manifestations: Report of three  
38 608 affected patients. *Am J Med Genet A* 2014; **164A**:1245–1253.
- 39  
40 609 50. Morimoto K, Shimizu T, Furukawa K, Morio H, Kurosawa H, Shirasawa T. Transgenic  
41 610 expression of the EXT2 gene in developing chondrocytes enhances the synthesis of heparan

- 1 sulfate and bone formation in mice. *Biochem Biophys Res Commun* 2002; **292**: 999–1009.
- 2 51. Salazar VS, Gamer LW, Rosen V. BMP signaling in skeletal development, disease and
- 3 repair. *Nat Rev Endocrinology* 2016; **12**:203–221.
- 4 52. Bragdon B, Moseychuk O, Saldanha S et al. Bone Morphogenetic Proteins: A critical review.
- 5 *Cell Signal* 2011; **23**: 609–620.
- 6 53. Ruvinsky A, Flood WD, Costantini F. Developmental mosaicism may explain spontaneous
- 7 reappearance of the AxinFu mutation in mice. *Genesis* 2001; **29**: 49-57.
- 8 54. Lu QL, Yao YH, Yao YF et al. Angiogenic Factor AGGF1 Promotes Therapeutic Angiogenesis
- 9 in a Mouse Limb Ischemia Model. *PLoS ONE* 2012; **7**: e46998.
- 10 55. González-Alonso J. Human thermoregulation and the cardiovascular system. *Exp Physiol*
- 11 2012; **97**:340–346.
- 12 56. Meadows JRS, Lindblad-Toh K. Dissecting evolution and disease using comparative
- 13 vertebrate genomics. *Nat Rev Genet* 2017; **18**, 624–636.
- 14 57. Dick DM, Bierut LJ. *Curr Psychiatry Rep* 2006; **8**: 151.
- 15 58. Horn KE, Glasgow SD, Gobert D et al. DCC expression by neurons regulates synaptic
- 16 plasticity in the adult brain. *Cell Rep* 2010; **31**:173-185.
- 17 59. Hori K, Hoshino M. Neuronal Migration and AUTS2 Syndrome. *Brain Sci* 2017; **7**:e54.
- 18 60. Hazen MM, Woodward AL, Hofmann I et al. Mutations of the hemophagocytic
- 19 lymphohistiocytosis-associated gene UNC13D in a patient with systemic juvenile idiopathic
- 20 arthritis. *Arthritis Rheum* 2008; **58**:567–570.
- 21 61. Procter M, Wolf B and Mao R. Forty-eight novel mutations causing biotinidase deficiency.
- 22 *Mol Genet Metab* 2016; **117**:369–372.
- 23 62. Munier S, Delcroix-Genete D, Carthagena L et al. Characterization of two candidate genes,
- 24 NCoA3 and IRF8, potentially involved in the control of HIV-1 latency. *Retrovirology* 2005;
- 25 **2**:73.
- 26 63. Fan ZY, Song YL. Chinese Primate Status and Primate Captive Breeding for Biomedical
- 27 Research in China. In: Institute for Laboratory Animal Research, National Research Council.
- 28 International Perspectives: The Future of Nonhuman Primate Resources. Washington DC:
- 29 National Academy Press. 2003.
- 30 64. Hao Xin. Monkey Research in China: Developing a Natural Resource. *Cell* 2007; **129**: 1033–
- 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

641 1036.

642 65. Zhou Q. Balancing the welfare: the use of non-human primates in research. *Trends Genet*

643 2014; **30**: 476–478.

644 66. Gradnigo JS, Majumdar A, Norgren RB Jr, Moriyama EN. Advantages of an Improved

645 Rhesus Macaque Genome for Evolutionary Analyses. *PLoS ONE* 2016; **11**: e0167376.

646 67. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform.

647 *Bioinformatics* 2009; **25**:1754–1760.

648 68. Li H, Handsaker B, Wysoker A et al. The sequence alignment/map format and SAMtools.

649 *Bioinformatics* 2009; **25**:2078–2079.

650 69. Mckenna A, Hanna M, Banks E, et al. The Genome Analysis Toolkit: a MapReduce

651 framework for analyzing next-generation DNA sequencing data. *Genome Res* 2010;

652 **20(9)**:1297–303.

653 70. Wang K, Li M, Hakonarson H. ANNOVAR: functional annotation of genetic variants from

654 high-throughput sequencing data. *Nucleic Acids Res* 2010; **38**: e164.

655 71. Patterson N, Price AL, Reich D. Population structure and eigenanalysis. *PLoS Genet* 2006;

656 **2**:e190.

657 72. Barrett JC, Fry B, Maller J et al. Haploview: analysis and visualization of LD and haplotype

658 maps. *Bioinformatics* 2005; **21**:263–265.

659 73. Ségurel L, Wyman M J, Przeworski M. Determinants of mutation rate variation in the human

660 germline. *Annu Rev Genomics Hum Genet* 2014; **15**: 47–70.

661 74. Weir BS, Cockerham CC. Estimating *F*-statistics for the analysis of population structure.

662 *Evolution* 1984; **38**:1358–1370.

663 75. Danecek P, Auton A, Abecasis G et al. The variant call format and VCFtools. *Bioinformatics*

664 2011; **27**:2156–2158.

665 76. Huang da W, Sherman BT, Lempicki RA. Systematic and integrative analysis of large gene

666 lists using DAVID Bioinformatics Resources. *Nat Protoc* 2009; **4**:44–57.

667 77. Liu Z, Tan X, Orozco-terWengel P, Zhou X, Zhang L, Tian S, et al. Supporting data for

668 “Population genomics of wild Chinese rhesus macaques reveals a dynamic demographic

669 history and local adaptation, with implications for biomedical research” *GigaScience*

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



671 **Tables**

672 Table1. Genetic diversity ( $\theta$ ) and effective population size ( $N_e$ ) in Chinese rhesus macaques  
 673 based on segregating sites ( $\theta_w$ ) and nucleotide diversity ( $\theta\pi$ ).  
 674

		Sample size (n)	$\theta_w$		$\theta\pi$	
			$\theta$	$N_e$	$\theta$	$N_e$
Chinese rhesus macaques (all samples)		81	0.00375	93,750	0.00247	61,750
	<i>M. m. littoralis</i>	29	0.00313	78,250	0.00240	60,000
	<i>M. m. tcheliensis</i>	5	0.00215	53,750	0.00230	57,500
subspecies	<i>M. m. brevicaudus</i>	5	0.00203	50,750	0.00207	51,750
	<i>M. m. lasiotis</i>	32	0.00298	74,500	0.00239	59,750
	<i>M. m. mulatta</i>	10	0.00303	75,750	0.00245	61,250

675

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

676 **Figure Legends**

677 **Figure 1.** Phylogeny and population genetic structure of 81 wild Chinese RMs. (a) Geographic  
678 distribution of RMs in China (gray shadow) and the 17 sampling sites along with their  
679 subspecies assignment. (b) Neighbor-joining (NJ) tree and clustering solution inferred using  
680 STRUCTURE and displaying five populations (inferred with Evanno's  $\Delta K$  method;  
681 Supplementary Fig. 5). (c) Principal component analysis plots depicting the first two components  
682 (variance explained by PC1 = 7.24% and PC2 = 5.69%).

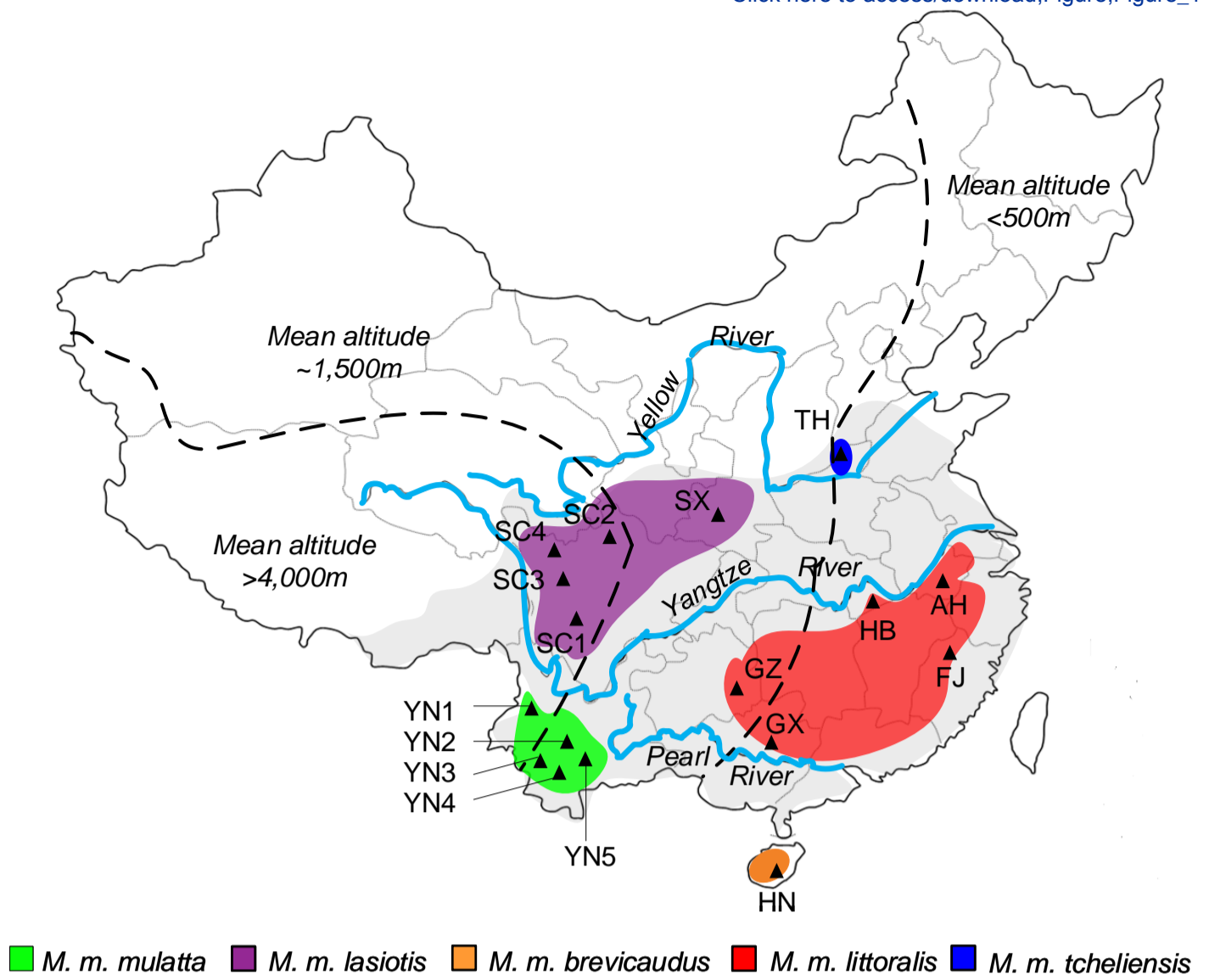
683 **Figure 2.** Demographic history and differentiation scenarios of Chinese RMs. (a) Historical  
684 changes in effective population size reconstructed using the pairwise sequential Markovian  
685 coalescent (PSMC) applied on individual whole genomes for each of the five subspecies. The  
686 generation length ( $g$ ) and the neutral mutation rate per generation ( $\mu$ ) were assumed to be 11  
687 years and  $1.08 \times 10^{-8}$ , respectively. The Xixiabangma Glaciation (XG, 1,200-800 kya),  
688 Penultimate Glaciation (PG, 200-130 kya) and Last Glaciation (LG, 70-10 kya) are shaded in  
689 gray. (b) Demographic history inferred by *fastsimcoal2*. The width of the gray bars and numbers  
690 on them indicate the estimated effective population size (all effective population sizes were  
691 converted to individuals). The arrows indicate migration rate between different subspecies. The  
692 detailed migration rates are listed in Supplementary Table 5. Numbers at the right show the  
693 divergence times between subspecies (all times were converted to years assuming a generation  
694 time of 11 years). (c) Biogeographic scenario for RMs. Chinese RMs separates from Indian RMs  
695 ~ 162 kya [13], followed by further migration into China by the different RM subspecies  
696 indicated with arrows colored following the color key in Fig. 1a.

697 **Figure 3.** Genomic regions with selection sweep signals in RM. (a) Distribution of  $\log_2(\theta_\pi M. m.$   
698 *lasiotis*/ $\theta_\pi M. m. tcheliensis$ ) and  $Z(F_{ST})$  of 50-kb windows with 25-kb steps. Blue dots located in  
699 the selected regions requirement (corresponding to  $Z$  test  $P < 0.05$ , where  $Z(F_{ST}) \geq 1.848$  and  $\theta_\pi$   
700 log-ratio  $\geq 1.203$ ) represent selected windows for *M. m. tcheliensis*. (b) Morphological  
701 comparison between *M. m. tcheliensis* and *M. m. lasiotis*. M and F represent males and females.  
702 (c) Example of genes with selection sweep signals. *Ext2*, *Rpgrip11*, *Fbp2* and *Fbp1* in *M. m.*  
703 *tcheliensis* and *Axin1*, *Aggf1* and *Hspa4* in *M. m. brevicaudus*.  $F_{ST}$  and  $\theta_\pi$  log-ratio between the  
704 two subspecies are represented in red and blue, respectively. All values in figure 3c are plotted  
705 using 50-kb windows with half steps. Genome annotations are show at the bottom (black bar,

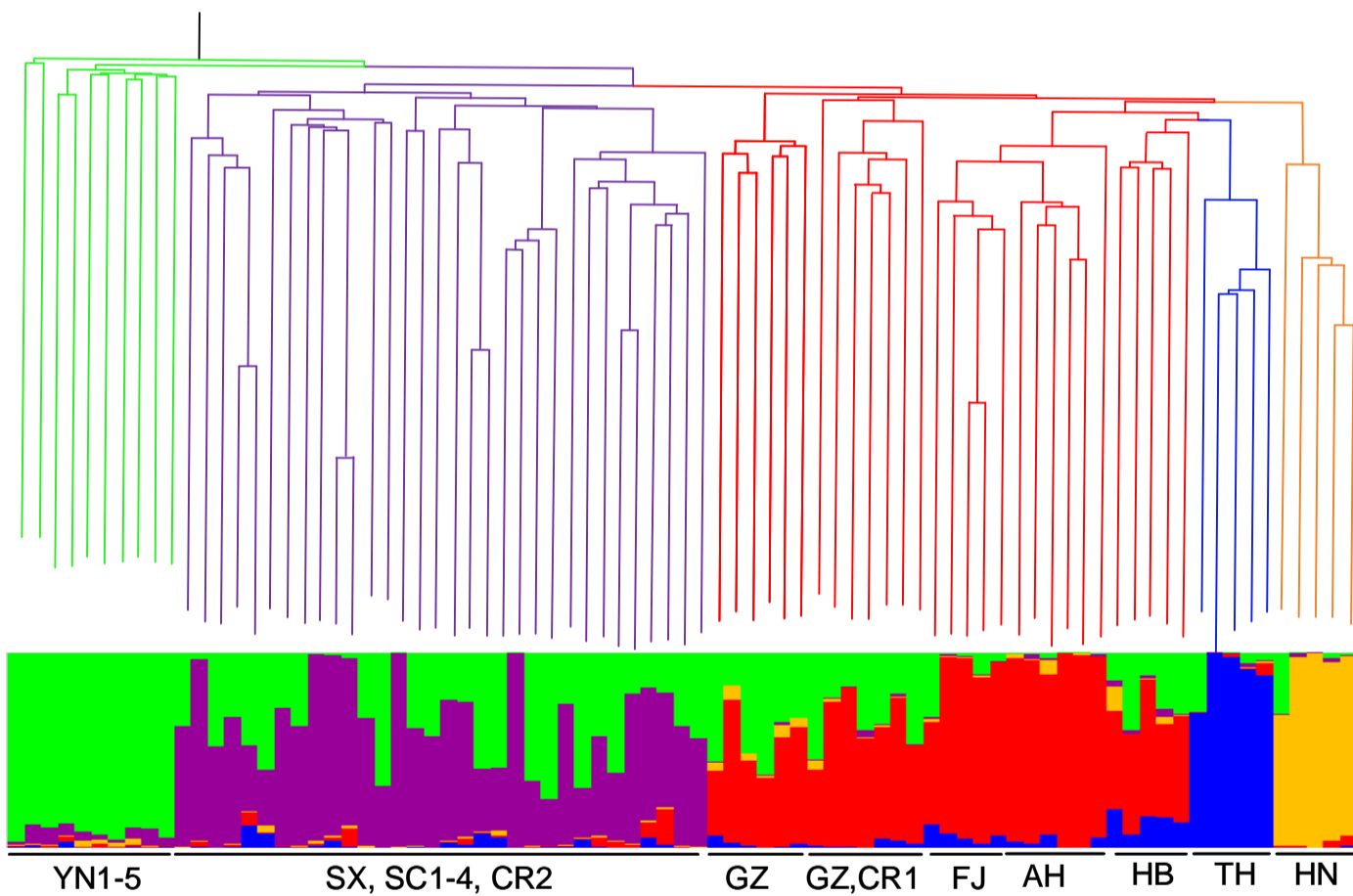
1 706 coding sequences (CDS); purple bar, genes). (d) SNP genotypes in putative selective sweeps  
2 707 containing *Ext2*, *Rpgrip11*, *Fbp2*, *Fbp1*, *Axin1*, *Aggf1* and *Hspa4*. (e) Non-synonymous variants  
3  
4 708 in gene *Ext2*, *Rpgrip11* and *Hspa4*.

5  
6 709 **Figure 4.** Population study of putative pathogenic SNPs in Chinese RM subspecies. (a) The site  
7  
8 710 and frequency of pathogenic SNPs located in *Unc13d* and *Btd* genes. (b) Scheme of the *Ncoa3*  
9  
10 711 gene in RM. The positions of nonsynonymous polymorphisms (black) and three amino-acid  
11  
12 712 deletions (in red) are marked. (c) Private and shared pathogenic SNPs in Chinese RM subspecies  
13  
14 713 (blue: *M. m. tcheliensis*; orange: *M. m. brevicaudus*; red: *M. m. littoralis*; green: *M. m. mulatta*;  
15  
16 714 purple: *M. m. lasiotis*). The sizes of the areas are not proportional to the magnitude of the  
17  
18 715 numbers. (d) NJ tree including the 81 Chinese RMs derived from this study, the 26 captive  
19  
20 716 Chinese RMs from Zhong et al. [7] are indicated by blue dot.

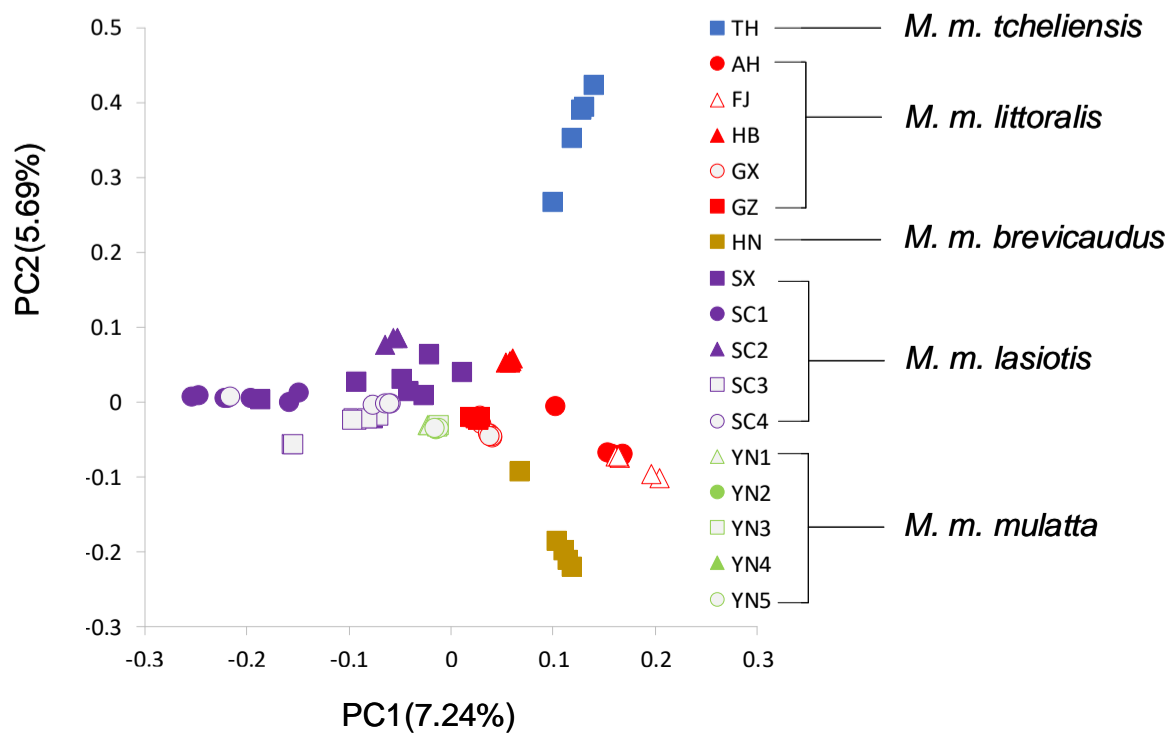
a

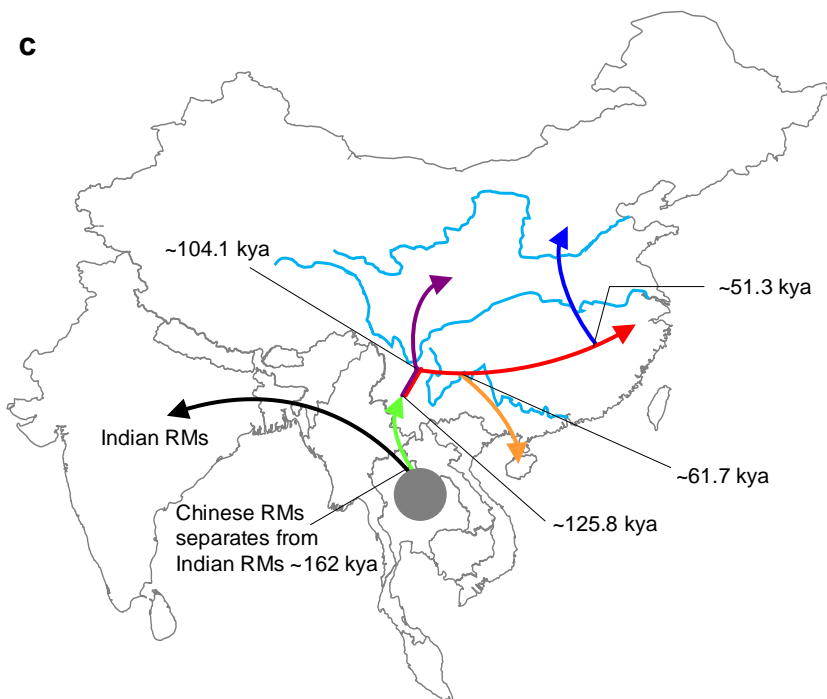
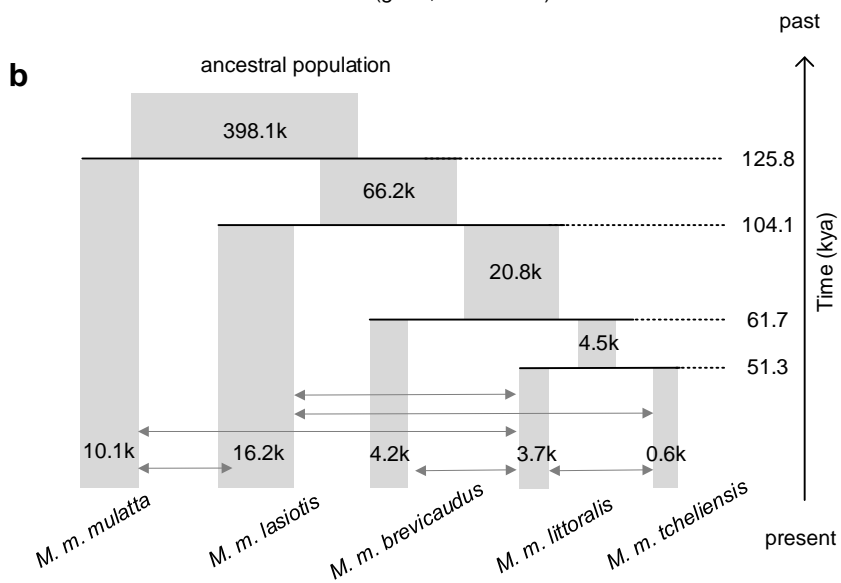
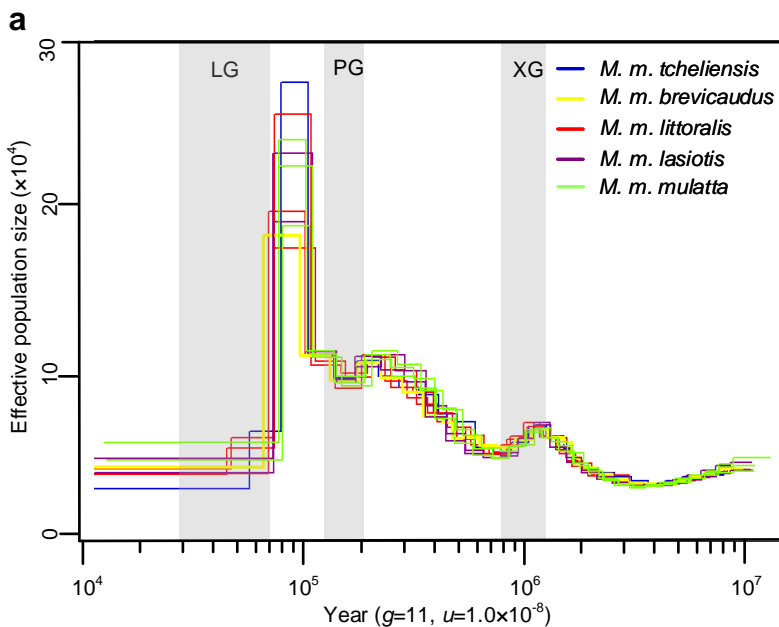


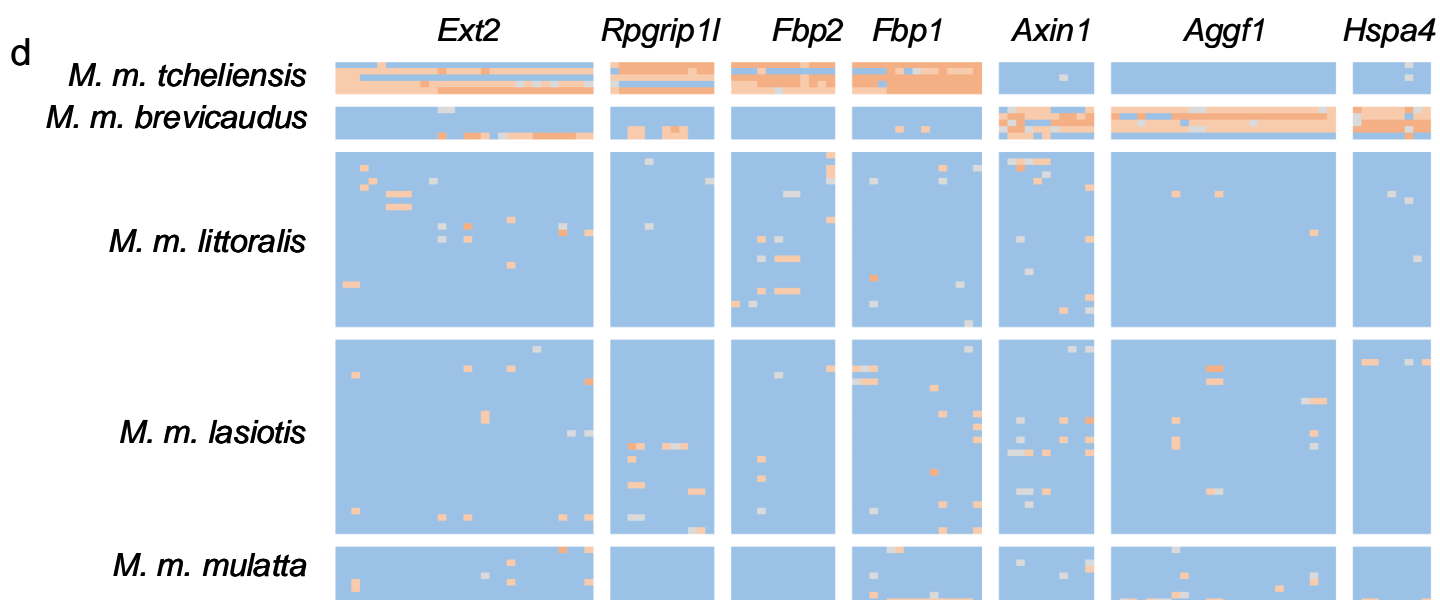
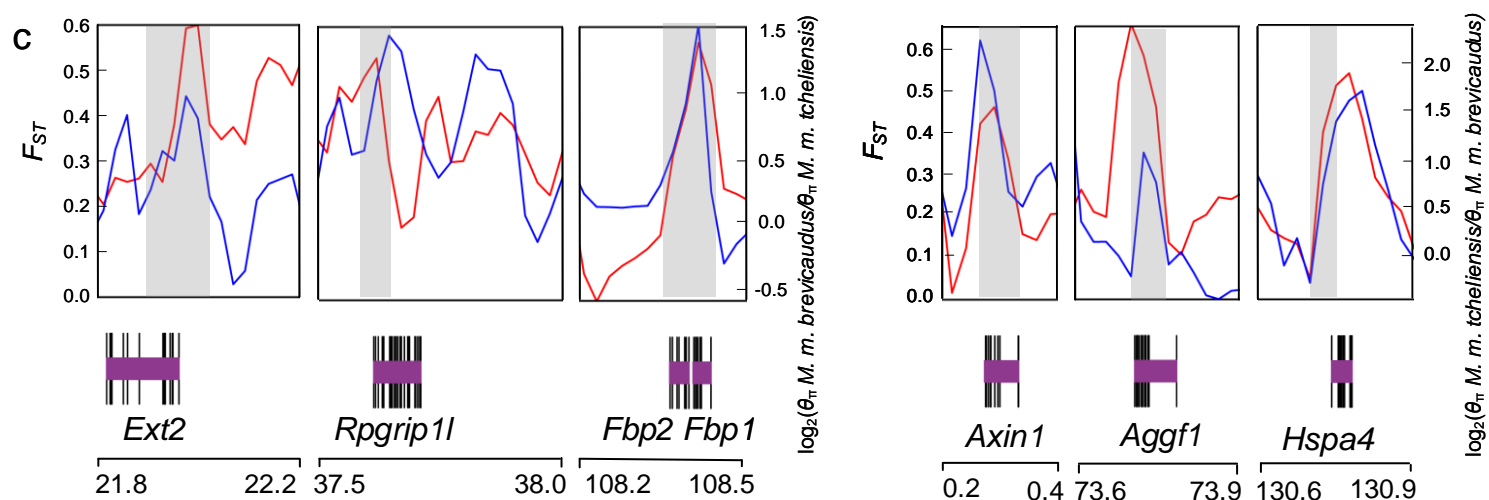
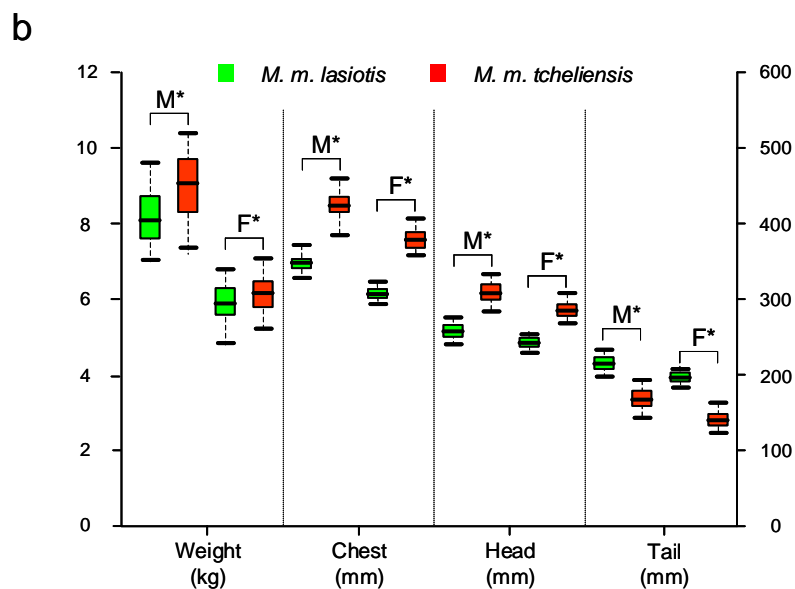
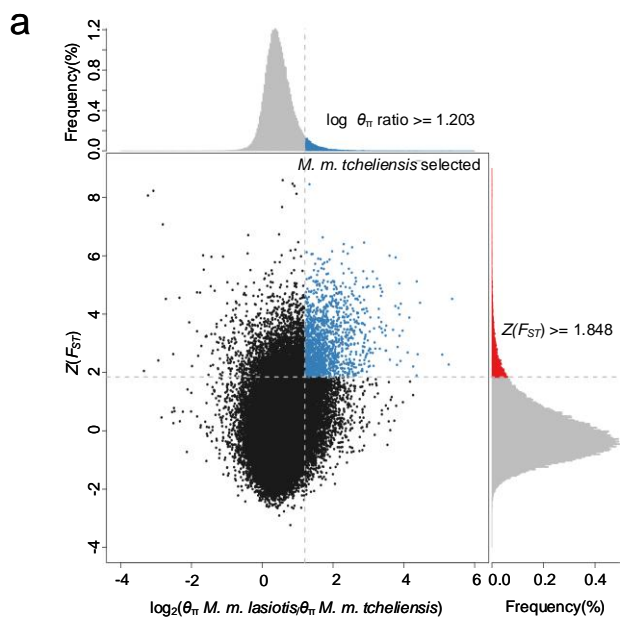
b



c



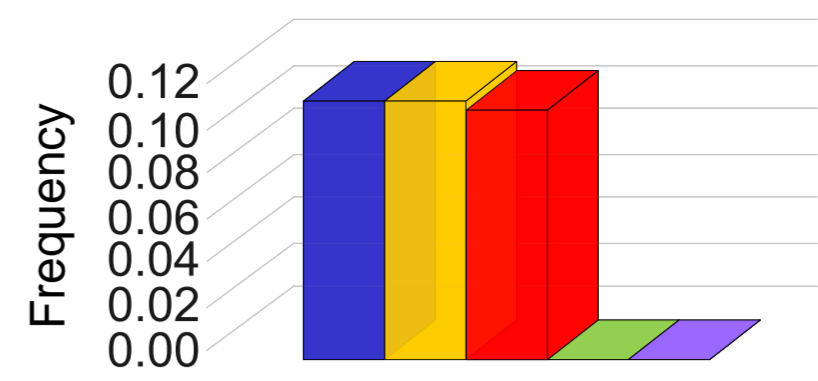




a

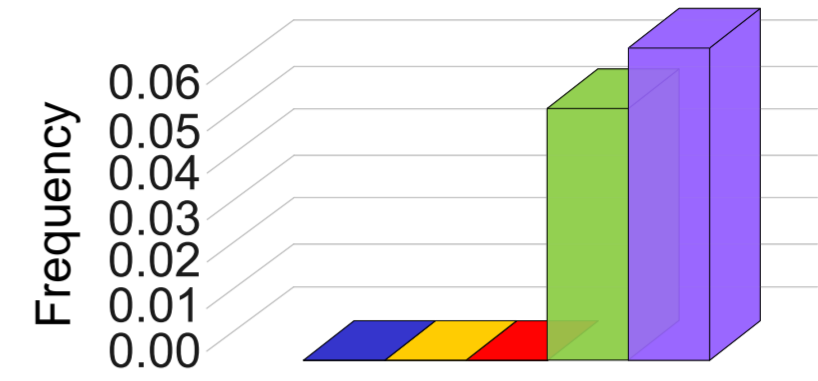
**UNC13D**

Human Reference **C A T C C T C C T C A C C T G C A G C C**  
 Human pathogenic SNP **. A . . . . . T . . . . .**  
 Macaque Reference **. C . . . . . C . . . . .**  
 Macaque pathogenic SNP **. C . . . . . T . . . . .**



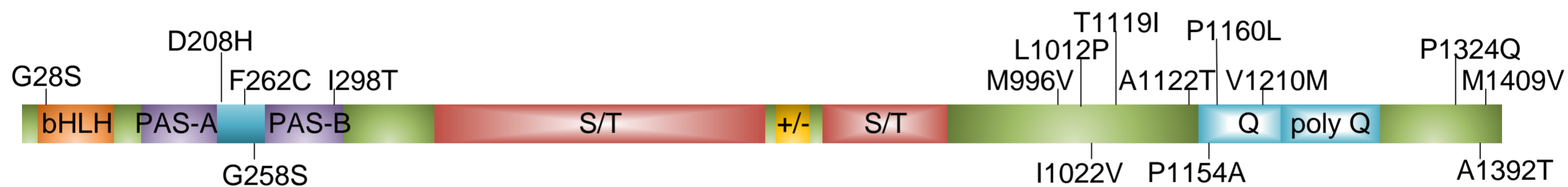
**BTD**

Human Reference **G G C A C T T A C T A C A T C C A A G T**  
 Human pathogenic SNP **. . . . . C . . . . . A . .**  
 Macaque Reference **. . . . . A . . . . . G . .**  
 Macaque pathogenic SNP **. . . . . C . . . . . G . .**

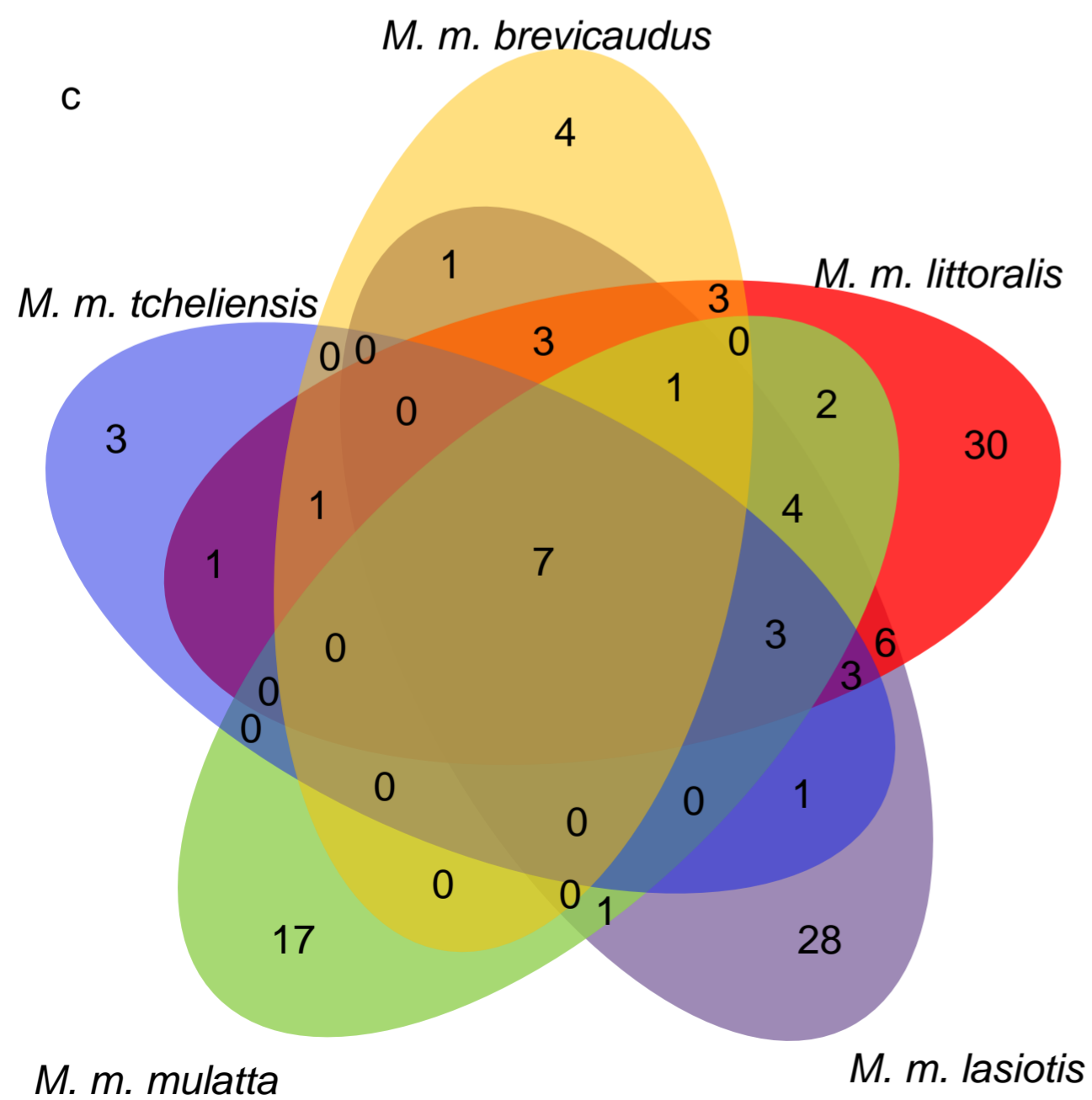


■ *M. m. tcheliensis*  
 ■ *M. m. brevicaudus*  
 ■ *M. m. littoralis*  
 ■ *M. m. mulatta*  
 ■ *M. m. lasiotis*

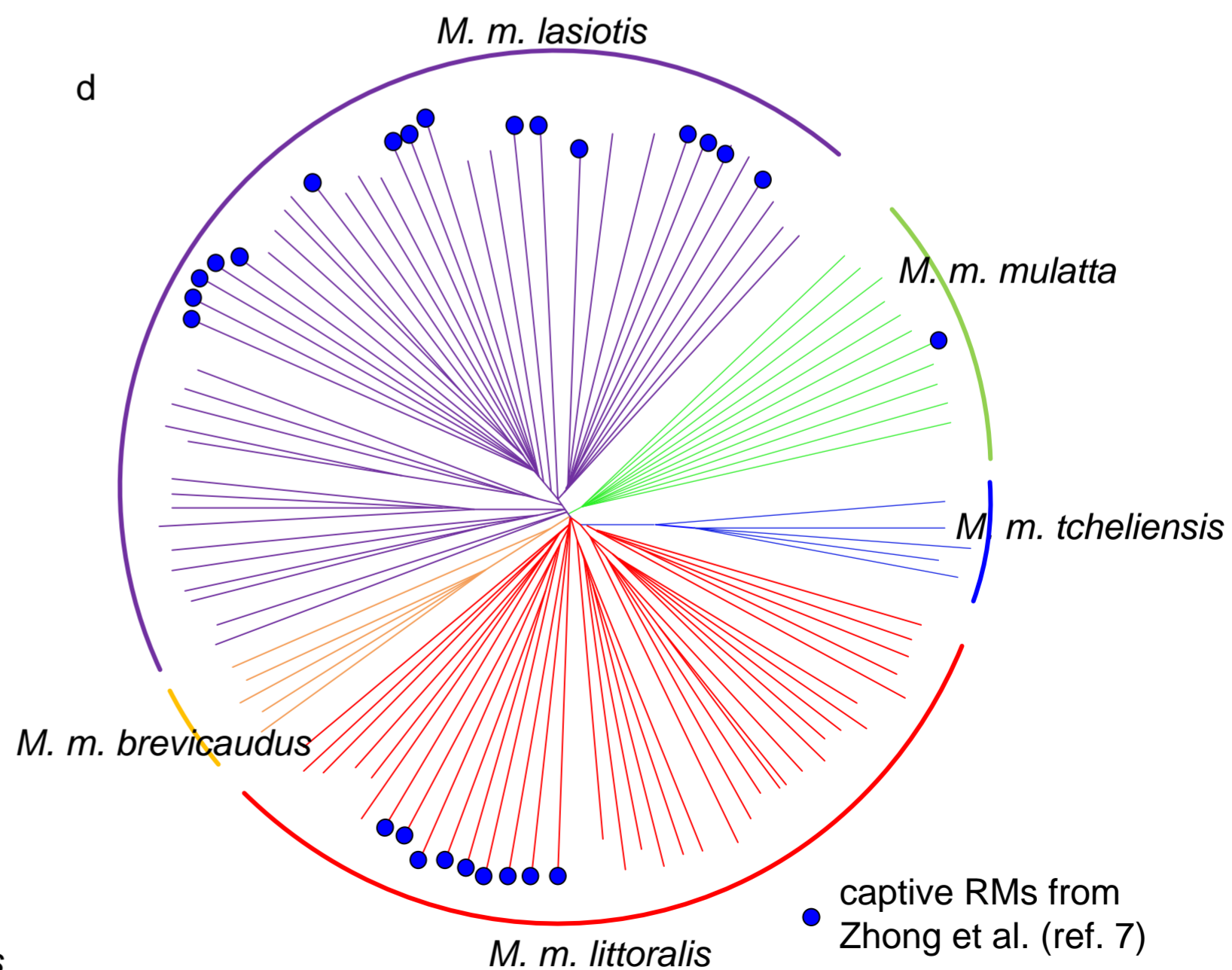
b




c

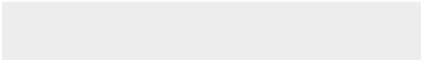



d





Click here to access/download  
**Supplementary Material**  
renamed\_cae38.docx







Click here to access/download  
**Supplementary Material**  
Supplemental Data 1.xlsx



Click here to access/download  
**Supplementary Material**  
Supplemental Data 2.xlsx