

# 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-17-00291	
<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>	Strategic Priority Research Program of the Chinese Academy of Sciences (XDPB0202)	Dr Ming Li
	NSFC (31530068)	Dr Ming Li
	NSFC (31471989)	PhD Zhijin Liu
	National Key R&D Program of China (2016YFC0503200)	Dr Ming Li
<b>Abstract:</b>	<p>The rhesus macaque (RM, <i>Macaca mulatta</i>) is the most important nonhuman primate model in evolutionary biology and biomedical research. We present the first population genomics survey of wild RMs, comprising 81 geo-referenced individuals representing five subspecies from 17 locations in China, covering a large fraction of the species' natural distribution. Genetic diversity, measured with a total of 55.4 M autosomal single nucleotide polymorphisms (SNPs), is higher in wild RMs than captive populations. We find a hierarchical population structure with four distinct genetic lineages found on the mainland and one on Hainan Island recapitulating current subspecies designations. The five subspecies are estimated to have diverged between 140 and 72 thousand years ago, but with recent gene flow among some groups. Consistent with the expectation of a larger body size in colder climates (Bergman's rule), the northernmost RM lineage (subspecies, <i>M. m. tcheliensis</i>) exhibits the largest body size of all Chinese RMs and was featured with positively selected genes responsible for skeletal development. The tropical subspecies <i>M. m. breviceaudus</i> was characterized by positively selected genes related to cardiovascular function and response to temperature stimuli, which are potentially involved in adaptation to tropical climates. We further delineated 111 RM SNPs matching human disease-causing variants with 74 being subspecies-specific. The data presented herein provides a reference resource for the choice of sub-group of RMs when carrying out biomedical experiments. The unexpected demographic history of Chinese RMs, coupled with their history of local adaptation offers new insights into the evolution of RMs and provides valuable baseline information for biomedical research.</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	

	Shilin Tian
	Liye Zhang
	Guangjian Liu
	Zhongze Yan
	Huailiang Xu
	Boshi Wang
	Baoping Ren
	Peng Zhang
	Zuofu Xiang
	Binghua Sun
	Christian Roos
	Michael W. Bruford
	Ming Li
<b>Order of Authors Secondary Information:</b>	
<b>Opposed Reviewers:</b>	
<b>Additional Information:</b>	
<b>Question</b>	<b>Response</b>
Are you submitting this manuscript to a special series or article collection?	No
<b>Experimental design and statistics</b>  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.  Have you included all the information requested in your manuscript?	Yes
<b>Resources</b>  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.  Have you included the information requested as detailed in our <a href="#">Minimum Standards Reporting Checklist</a> ?	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 a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p>	<p>Yes</p>

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, 3, †</sup>, Pablo Orozco-terWengel<sup>4, †</sup>, Xuming Zhou<sup>5, †</sup>, Shilin Tian<sup>6, 7, †</sup>, Liye  
6 Zhang<sup>1, 2</sup>, Guangjian Liu<sup>6</sup>, Zhongze Yan<sup>1, 3</sup>, Huailiang Xu<sup>7</sup>, Boshi Wang<sup>1</sup>, Baoping Ren<sup>1</sup>, Peng  
7 Zhang<sup>8</sup>, Zuofu Xiang<sup>9</sup>, Binghua Sun<sup>10</sup>, Christian Roos<sup>11, \*</sup>, Michael W. Bruford<sup>4, \*</sup>, Ming Li<sup>1, \*</sup>

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

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

11 <sup>3</sup> Institute of Health Sciences, Anhui University, Hefei, 230601, China.

12 <sup>4</sup> School of Biosciences, Cardiff University, Sir Martin Evans Building, Museum Avenue, Cardiff  
13 CF10 3AX, United Kingdom.

14 <sup>5</sup> Division of Genetics, Department of Medicine, Brigham and Women's Hospital, Harvard Medical  
15 School, Boston, MA 02115, USA.

16 <sup>6</sup> Novogene Bioinformatics Institute, Beijing 100083, China.

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

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

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

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

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

24 <sup>†</sup> Contributed equally

25 \* Correspondence: Ming Li, [lim@ioz.ac.cn](mailto:lim@ioz.ac.cn); Michael W. Bruford, [BrufordMW@cardiff.ac.uk](mailto:BrufordMW@cardiff.ac.uk);  
26 Christian Roos, [CRoos@dpz.eu](mailto:CRoos@dpz.eu)

1       27   **Abstract**

2  
3       28   The rhesus macaque (RM, *Macaca mulatta*) is the most important nonhuman primate model in  
4  
5       29   evolutionary biology and biomedical research. We present the first population genomics survey of  
6  
7       30   wild RMs, comprising 81 geo-referenced individuals representing five subspecies from 17 locations  
8  
9       31   in China, covering a large fraction of the species' natural distribution. Genetic diversity, measured  
10  
11      32   with a total of 55.4 M autosomal single nucleotide polymorphisms (SNPs), is higher in wild RMs  
12  
13      33   than captive populations. We find a hierarchical population structure with four distinct genetic  
14  
15      34   lineages found on the mainland and one on Hainan Island recapitulating current subspecies  
16  
17      35   designations. The five subspecies are estimated to have diverged between 140 and 72 thousand years  
18  
19      36   ago, but with recent gene flow among some groups. Consistent with the expectation of a larger body  
20  
21      37   size in colder climates (Bergman's rule), the northernmost RM lineage (subspecies, *M. m. tcheliensis*)  
22  
23      38   exhibits the largest body size of all Chinese RMs and was featured with positively selected genes  
24  
25      39   responsible for skeletal development. The tropical subspecies *M. m. breviceaudus* was characterized  
26  
27      40   by positively selected genes related to cardiovascular function and response to temperature stimuli,  
28  
29      41   which are potentially involved in adaptation to tropical climates. We further delineated 111 RM  
30  
31      42   SNPs matching human disease-causing variants with 74 being subspecies-specific. The data  
32  
33      43   presented herein provides a reference resource for the choice of sub-group of RMs when carrying  
34  
35      44   out biomedical experiments. The unexpected demographic history of Chinese RMs, coupled with  
36  
37      45   their history of local adaption offers new insights into the evolution of RMs and provides valuable  
38  
39      46   baseline information for biomedical research.

40  
41  
42      47   **Keywords:** *Macaca mulatta*, population genomics, adaptive selection, biomedical model  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

## 48 **Introduction**

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

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

71 To date, the genomes of 133 captive RMs from eight colonies have been sequenced, however,  
72 124 of them are of Indian-origin and only nine individuals were presumed to be of Chinese origin  
73 [6]. Recently, Zhong *et al.* [7] reported genomic variation in 26 Chinese captive RMs identifying  
74 ~46 M (million) single nucleotide polymorphisms (SNPs). Nevertheless, most of these RM genetic  
75 variation is limited to captive populations which may contain composite genotypes due of admixture  
76 among animals of different and unclear origin [20]. Here we present the first attempt to survey the

1 77 geo-referenced genomic diversity in wild Chinese RM populations, which is the largest extant  
2 78 population of the species. The current effective population size of Chinese and Indian RM was  
3  
4 79 estimated to be approximate 240,000 and 17,000 individuals, respectively, indicating that the  
5  
6 80 Chinese RMs are likely to harbor substantially more genomic diversity compared to their Indian  
7  
8 81 conspecifics [13]. Therefore, this population genomic survey of 81 RMs originating from 17 wild  
9  
10 82 locations across China including phylogenetic, demographic and genome-wide selection scans,  
11  
12 83 corresponds to the most comprehensive characterization of RM genetic diversity to date and aimed  
13  
14 84 at characterizing the processes leading to the extant patterns of variability, as well as identifying the  
15  
16 85 potential implications for the use of these populations in biomedical research.  
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

## 86 Results and Discussion

### 87 Genetic diversity, phylogeny and population structure

88 Blood and tissue samples from 79 wild-born RMs, representing five subspecies [21, 22], were  
89 collected at 17 sites in China (*M. m. tcheliensis*: TH; *M. m. littoralis*: AH, FJ, HB, GX, GZ; *M. m.*  
90 *brevicaudus*: HN; *M. m. lasiotis*: SX, SC1, SC2, SC3, SC4; *M. m. mulatta*: YN1, YN2, YN3, YN4,  
91 YN5; Fig. 1a). Genome sequences of two additional Chinese RMs (CR1 and CR2) were retrieved  
92 from NCBI [9, 23, 24]. Re-sequencing was at a high average depth of  $28.26 \pm 4.75 \times$  for ten  
93 individuals and a moderate average depth of  $10.64 \pm 1.16 \times$  for the remainder ( $n=71$ ), with an overall  
94 average genome coverage of 95.43% of the RM reference (rheMac2, Supplementary Table 1). A  
95 total of 55,404,179 SNPs were identified and used for downstream analyses (Supplementary Table  
96 2), with wild RMs carrying on average over 20% more genomic variation than captive individuals  
97 (43.7 M SNPs detected in 31 Indian and Chinese captive RMs and 46.1 M SNPs found in 133 Indian  
98 and Chinese captive RMs)[6,7]. The number of SNPs per individual ranged from 7.3 to 11.6 M  
99 (mean of 9.5 M; Supplementary Fig. 1 and Supplementary Table 3). Among Chinese RM subspecies,  
100 *M. m. mulatta* had the highest heterozygosity ( $0.202\% \pm 1.09 \times 10^{-4}$ ), followed by *M. m. littoralis*  
101 ( $0.180\% \pm 1.55 \times 10^{-4}$ ) and *M. m. lasiotis* ( $0.178\% \pm 1.47 \times 10^{-4}$ ). The lowest heterozygosity rates were  
102 found in *M. m. brevicaudus* ( $0.157\% \pm 1.17 \times 10^{-4}$ ) and *M. m. tcheliensis* ( $0.135\% \pm 3.07 \times 10^{-4}$ )  
103 (Supplementary Fig. 2). Among all detected SNPs, 7,575,099 were shared among all subspecies and  
104 23,676,191 were shared by at least two subspecies, with the remaining SNPs confined to a single  
105 subspecies (Supplementary Fig. 3a). For each subspecies, the subspecies-specific SNPs (ssSNPs)  
106 ranged from 834,655 to 8,507,232 and the non-synonymous ssSNPs varied from 3,723 to 27,537  
107 (Supplementary Fig. 3a, b).

108 We reconstructed a neighbor-joining (NJ) tree for Chinese RMs based on autosomal SNPs,  
109 using Indian RMs and *M. sylvanus* as outgroups (Fig. 1b and Supplementary Fig. 4). Individuals  
110 from *M. m. lasiotis*, *M. m. brevicaudus* and *M. m. tcheliensis* form monophyletic lineages  
111 respectively, while *M. m. mulatta* and *M. m. littoralis* are paraphyletic. The divergence among  
112 Chinese RMs started with the successive splitting of the *M. m. mulatta* lineages, followed by *M. m.*  
113 *lasiotis*, before the eastern subspecies differentiated. Among the latter, *M. m. brevicaudus* and *M.*  
114 *m. tcheliensis* diverged from *M. m. littoralis*, respectively. Next, we performed a population



1 115 structure analysis using STRUCTURE (version 2.3.4) [25], which estimates individual ancestry and  
2  
3 116 admixture structure analysis proportions assuming  $K$  ancestral populations. Plots of  $\Delta K$  generated  
4  
5 117 from STRUCTURE results indicated five genetic clusters present in the full data set (Fig. 1b and  
6  
7 118 Supplementary Fig. 5). A principal component analysis (PCA) corroborated the division of Chinese  
8  
9 119 RMs into five groups. The first eigenvector separated *M. m. mulatta* and *M. m. lasiotis* from *M. m.*  
10  
11 120 *tcheliensis*, *M. m. littoralis* and *M. m. brevicaudus* (variance explained = 6.40%, Tracy-Widom  $P =$   
12  
13 121  $4.29 \times 10^{-42}$ ), and the second eigenvector further separated *M. m. tcheliensis*, *M. m. littoralis* and *M.*  
14  
15 122 *m. brevicaudus* (variance explained = 5.07%, Tracy-Widom  $P = 5.29 \times 10^{-22}$ ) (Fig. 1c,  
16  
17 123 Supplementary Table 4). The division of Chinese RMs into five geographic lineages supports the  
18  
19 124 former taxonomic division of Chinese RMs into five subspecies [21, 22]. *M. m. mulatta* (YN1-5)  
20  
21 125 and *M. m. lasiotis* (SC1-4, SX) form the pan-western populations of Chinese RMs, with both  
22  
23 126 subspecies inhabiting the montane Tibetan Plateau regions with an altitude  $\geq 1500$  meters above sea  
24  
25 127 level in western China and separated from each other by the Yangtze River. *M. m. littoralis* (AH,  
26  
27 128 FJ, HB, GX, GZ), *M. m. tcheliensis* (TH) and *M. m. brevicaudus* (HN) occur in the eastern coastal  
28  
29 129 lowland of China and form the pan-eastern population. *M. m. tcheliensis* from the Taihang  
30  
31 130 Mountains area (TH) is the northernmost ( $34^{\circ}54' - 35^{\circ}16' \text{ N}$ ;  $112^{\circ}02' - 112^{\circ}52' \text{ E}$ ), while *M. m.*  
32  
33 131 *brevicaudus*, restricted to Hainan Island, is the most southern Chinese RM subspecies.

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

### 133 **Demographic and phylogeographic history**

134 We applied the pairwise sequential Markovian coalescent (PSMC) [26] using ten RM individuals  
135 with an average sequencing coverage depth higher than  $20\times$  (one individual of *M. m. tcheliensis* and  
136 one of *M. m. brevicaudus*, two of *M. m. lasiotis*, three of *M. m. littoralis* as well as three individuals  
137 of *M. m. mulatta*) to infer the ancient demographic history of Chinese RMs. The inferred PSMC  
138 trajectories were very similar for all analyzed individuals throughout most of the species' history  
139 until  $\sim 110$  kya reflecting the species' cohesiveness (Fig. 2a). The ancient demographic history of  
140 RMs is marked by population fluctuations following the glacial periods during the Pleistocene [27].  
141 Approximately 800-500 kya all Chinese RMs experienced a population reduction at the time of the  
142 Naynayxungla Glaciation (NG, 780-500 kya), followed by an expansion during the Mid-Pleistocene  
143 inter-glaciation (500-200 kya). This expansion was then interrupted by the Penultimate Glaciation  
144 (PG, 200-130 kya) when suitable habitat might have been lost leading to a population decline [27].

1 145 PSMC analyses also suggested that while *M. m. mulatta* and *M. m. lasiotis* stabilized with effective  
2 146 population sizes somewhere around 110 kya, *M. m. tcheliensis*, *M. m. littoralis* and *M. m.*  
3  
4 147 *brevicaudus* went through a dramatic population increase and a subsequent bottleneck reaching a  
5  
6 148 stable effective population size somewhere around 50-60 kya (Fig. 2a).

7  
8 149 Due to the observed divergence of the five Chinese RM subspecies, we further employed the  
9  
10 150 joint site frequency spectrum (SFS) approach to model scenarios that could explain the observed  
11  
12 151 population structure, as well as respective divergence times between the five RM lineages. These  
13  
14 152 analyses, carried out using *fastsimcoal2* [28] produced a significantly better fit of a step by step  
15  
16 153 divergence scenario than alternative ones (Supplementary Tables 5 and 6, Supplementary Fig. 6),  
17  
18 154 and support the demographic changes observed with the PSMC analyses. Under this model,  
19  
20 155 following the divergence between the ancestral lineages of the Indian and Chinese RM (~162 kya),  
21  
22 156 the ancestor of the *M. m. mulatta* lineage diverged from that of the remaining Chinese RMs ~140.2  
23  
24 157 kya near the end of the last interglacial (Fig. 2b) [6, 13]. This divergence was associated with a mild  
25  
26 158 decrease in effective population size ( $N_e$ ) of *M. m. mulatta* with respect to its ancestor. Subsequently,  
27  
28 159 *M. m. lasiotis* diverged from the ancestral lineage of pan-eastern RM ~107.1 kya, undergoing a  
29  
30 160 similar reduction in  $N_e$  as *M. m. mullata*. Contrastingly, this divergence was followed by an almost  
31  
32 161 doubling of the ancestral effective population size ( $N_{A2}$  = 232.5k) of the ancestral pan-eastern RM  
33  
34 162 lineage (Fig 2b). The divergence time among *M. m. tcheliensis*, *M. m. littoralis* and *M. m.*  
35  
36 163 *brevicaudus* was estimated to occurred ~71.7 kya, at the start of the period leading to the last glacial  
37  
38 164 maximum (fig 2b) [29, 30], and coinciding with a drastic decrease in  $N_e$  in these three lineages.  
39  
40 165 Gene flow after the divergence of subspecies occurred among almost all five lineages.

41  
42 166 A previous study of mitochondrial DNA identified two major haplogroups dividing Chinese  
43  
44 167 RMs in a western and an eastern clade, and with modern Chinese RMs thought to have undergone  
45  
46 168 a northward expansion while entering China via two possible routes: the first into the western  
47  
48 169 mountains and the second following the eastern coast [31]. Our evolutionary model, however,  
49  
50 170 suggests a “step-by-step” colonization process of RMs into China (Fig 2c). After the divergence  
51  
52 171 from the Indian population (~162 kya) [6, 13], the ancestor of Chinese RMs colonized the Tibetan  
53  
54 172 Plateau from southwestern China, and then experienced a range expansion north and eastwards. The  
55  
56 173 pan-western population (*M. m. mulatta* and *M. m. lasiotis*) inhabited the western montane region in  
57  
58 174 China, while the pan-eastern population (*M. m. tcheliensis*, *M. m. littoralis* and *M. m. brevicaudus*)

1 175 entered the eastern coastal region. Barriers such as the Yellow, Yangtze and Pearl rivers and open  
2 176 sea (Fig. 1a) led to further differentiation, limiting gene flow among them. Water bodies and  
3  
4 177 mountains could therefore be described as driving the formation of a habitat ‘lattice’ with the  
5  
6 178 different subspecies of RMs occupying different grids in the lattice.  
7

8  
9 179

## 10 180 **Signatures of selection and local adaptation**

11 181 The wide distribution of Chinese RMs and their respective contrasting habitat types, as well as their  
12  
13 182 wide use in biomedical studies, makes them an important case study for the analysis of signatures  
14  
15 183 of local adaptation to divergent selective pressures [32-34]. We identified putative targets of  
16  
17 184 selection by carrying out pair-wise comparisons between RM subspecies inhabiting the most  
18  
19 185 different environments to increase the chance of finding selection signatures, i.e., *M. m. tcheliensis*  
20  
21 186 that lives in the northernmost range of the species in the coldest environments, and *M. m.*  
22  
23 187 *brevicaudus* that inhabits the southernmost range of the species in a tropical environment. For each  
24  
25 188 analysis, we compared the five subspecies using the fixation index ( $F_{ST}$ ) and genetic diversity ( $\theta_{\pi}$ ),  
26  
27 189 calculated on 50kb long sliding windows (Fig. 3 and Supplementary Figs. 7-12). The top 5% of the  
28  
29 190 windows with the largest  $F_{ST}$  and  $\theta_{\pi}$  ratios ( $\theta_{\pi 2} / \theta_{\pi 1}$ ) in each pair-wise comparison were considered  
30  
31 191 to be potentially under positive selection. For each subspecies, we identified the intersection of  
32  
33 192 potential selective-sweep regions generated by all the pair-wise comparisons between a subspecies  
34  
35 193 and each of the other subspecies (four pairwise comparisons in each case) (Supplementary Fig. 7).  
36  
37 194 We used these consistent selective-sweep regions for further analyses, as they represent robust  
38  
39 195 putative positively selected regions. The sizes of candidate selective-sweep regions ranged from  
40  
41 196 0.275 Mb to 8.575 Mb and the number of genes located in these regions, which are expected to  
42  
43 197 represent targets of selection for each subspecies, varied from 8 to 141 in different subspecies  
44  
45 198 (Supplementary Table 7).  
46  
47  
48

49 199 *M. m. tcheliensis* from the Taihang (TH) Mountains area is the northernmost population of the  
50  
51 200 species. The TH mountain are characterized by a continental monsoon climate, and conditions for  
52  
53 201 RMs are harsh during winter and early spring with average temperatures of  $-20^{\circ}\text{C}$ . According to  
54  
55 202 Bergman’s rule, animals living in cold climates tend to have larger body sizes compared to their  
56  
57 203 relatives in warm climates (i.e. they have a lower surface area to volume ratio), so they radiate less  
58  
59 204 body heat per unit of mass [35]. Consistent with this expectation, of all RM subspecies, *M. m.*

1 205 *tcheliensis* exhibits the largest body size and mass, the shortest tail length, the longest forearm length  
2 206 and the largest head and chest circumference (Fig. 3b and Supplementary Table 8) [36,37]. Among  
3 207 the consistent signatures of positive selection identified in *M. m. tcheliensis* (128 genes), we found  
4 208 signatures of selective sweeps in 14 genes linked to limb morphogenesis and bone development  
5 209 (Supplementary Table 7), which present two highly enriched functional categories, “embryonic  
6 210 hind-limb morphogenesis” (three genes, modified Fisher Exact  $P=1.59\times 10^{-2}$ ) and i.e. “bone  
7 211 development” (three genes, modified Fisher Exact  $P=3.40\times 10^{-2}$ ) (Supplementary Table 9). Among  
8 212 these genes, *Papss2* is known to affect the development of the skeletal system in mouse and human  
9 213 and *Papss2* mutations could cause brachyolmia [38,39], while *Sox5* (Fig. 3c, d) plays an essential  
10 214 role in synovial joint morphogenesis via promoting both growth plate and articular chondrocyte  
11 215 differentiation [40]. *Bcl2* has been shown to regulate chondrocyte maturation during skeletal  
12 216 development and could influence long bone length [41]. These genes involved in the growth and  
13 217 development of the skeletal system and appendages are likely contributors to the larger body size of  
14 218 *M. m. tcheliensis*, and represent an undescribed adaptive pathway for primates living in colder  
15 219 climates.

16 220 In contrast, *M. m. brevicaudus* inhabits the tropical island of Hainan (HN) where it copes with a  
17 221 mean annual temperature of 24°C. We found 141 putatively selected genes in *M. m. brevicaudus*  
18 222 (Fig. 3c, d, Supplementary Table 7), seven of which were found in gene ontology (GO) terms related  
19 223 to cardiovascular system and blood circulation. For example, *Ppp3cb* related to GO term “heart  
20 224 development” and *Cttna3* related to GO term “regulation of heart rate by cardiac conduction” [42].  
21 225 In addition, *Camk2g* and *Ero1a* are directly involved in the GO terms “regulation of cellular  
22 226 response to heat” and “response to temperature stimulus”. We thus hypothesize that the  
23 227 cardiovascular system of *M. m. brevicaudus* might play an important role in stabilizing body  
24 228 temperature, assisted by blood flow through different body parts requiring good fluidity and vascular  
25 229 permeability to transfer heat out of the body [43]. Test of these hypothesis needs further functional  
26 230 assays, however, these genes, together with the positively selected genes identified in *M. m.*  
27 231 *tcheliensis*, are known to be relevant to human physical function, and thus are likely of importance  
28 232 in the adaptation of Chinese RMs to different climates.

29 233 Besides the genes related to the adaptation to various climate conditions, we also found positive  
30 234 selection in genes related to the nervous system. In *M. m. littoralis* three of the 104 identified

1 235 candidate genes are enriched in GO term “regulation of synaptic plasticity” (modified Fisher Exact  
2 236  $P=1.38E-02$ ; Supplementary Table 10) and four genes are enriched in the KEGG pathway  
3  
4 237 “serotonergic synapse” (modified Fisher Exact  $P=2.11E-02$ ; Supplementary Table 10). In *M. m.*  
5  
6 238 *tcheliensis* twelve putatively selected genes (Supplementary Table 7) are involved in the process of  
7  
8 239 neuron morphogenesis and synaptic transmission, and one of these gene, *Clstm2* is  
9  
10 240 the synaptic protein and reported to play an important role in learning and memory [44]. For *M. m.*  
11  
12 241 *brevicaudus*, seven putatively selected genes related to nervous system development were found.  
13  
14 242 For example, *Dcc* is reported to be required for long-term potentiation and memory [45]. *Auts2*, one  
15  
16 243 of the eight putatively selected genes in *M. m. mulatta*, has been shown to regulate neuronal  
17  
18 244 migration, and mutations in this gene cause mental dysfunction in human [46] (Supplementary Table  
19  
20 245 7). Our findings suggest that RM subspecies have experienced different adaptive processes in the  
21  
22 246 nervous system and respective genomic differences should be taken into account when animals are  
23  
24 247 selected for neurobiological research.  
25  
26  
27  
28

29

#### 30 **Disease-causing variants and implication for biomedical research**

31 250 Given the large evolutionary similarity between macaques and humans, human diseases are better  
32  
33 251 modeled in RMs than in many other animals. Thus, variants in RMs that match to orthologous  
34  
35 252 human variants annotated as ‘pathogenic’ are of particular interest. We examined presumed  
36  
37 253 homologous Chinese RM SNPs in the human genome and a total of 32,845,501 RM SNPs analyzed  
38  
39 254 in this study were successfully identified in the human genome (hg19). Among these SNPs, 111  
40  
41 255 variants matched human variants with the accordant reference alleles and alternative alleles were  
42  
43 256 annotated as ‘disease causing’ in HGMD or pathogenic in ClinVar. Those 111 RM SNPs affect genes  
44  
45 257 that cause specific human diseases including acromesomelic dysplasia maroteaux type, anonychia,  
46  
47 258 atransferrinemia, blau syndrome, Carcinoma of colon, Charcot-Marie-Tooth disease, deafness, early  
48  
49 259 infantile epileptic encephalopathy 7, glycogen storage disease and others (Supplementary Table 11).  
50  
51 260 Thirty-nine out of these 111 SNPs were identified in previous studies [6], while the remaining 72  
52  
53 261 SNPs are newly described here. Only seven pathogenic SNPs are shared by all five subspecies,  
54  
55 262 while 74 are subspecies-specific (Fig. 4c, Supplementary Table 11). For example, the SNP  
56  
57 263 rs116229331 in the *Unc13d* gene (human Chr17: 73836585C>T), known to cause juvenile  
58  
59 264 idiopathic arthritis in humans [47], has a RM homologue (RM Chr16: 71160253C>T, Fig. 4a) that  
60  
61  
62  
63  
64  
65

1 265 is present in *M. m. tcheliensis*, *M. m. brevicaudus* and *M. m. littoralis*, but absent in *M. m. lasiotis*  
2 266 and *M. m. mulatta*. Another pathogenic variant (rs397514345, human Chr3: 15686724 A>C) in the  
3  
4 267 *Btd* gene is involved in biotinidase deficiency [48]. Its homologous RM variant (RM Chr2:  
5  
6 268 157981062 A>C, Fig. 4a) is found only in *M. m. lasiotis* and *M. m. mulatta*. In addition, we also  
7  
8 269 identified 16 non-synonymous SNPs in the *Noca3* gene, which encodes a protein that modulates the  
9  
10 270 replication and transcriptional reactivation of HIV-1 during virus latency [49] (Fig. 4b). Ten of these  
11  
12 271 16 non-synonymous SNPs are private to one subspecies (Supplementary Table 12). The effects of  
13  
14 272 these variants on HIV-1 replication and reactivation are unknown and need further investigation, but  
15  
16 273 the high number of mutations suggests a complex response of the host to the virus.

17  
18 274 Overall, these findings suggest that the genomic architecture of Chinese RMs used in  
19  
20 275 biomedical research and their geographic origin could strongly influence the outcome of biomedical  
21  
22 276 experiments and should be taken into account when using Chinese RMs in clinical and  
23  
24 277 neurobiological research. Unfortunately, genome wide screening of RMs used in biomedical  
25  
26 278 research is so far only rarely conducted and uncharacterized animals are most often used.  
27  
28 279 Importantly, individuals from all five Chinese RM subspecies are used in biomedical research [50,  
29  
30 280 51]. Combined with our data, nine of the 26 captive Chinese RMs reported by Zhong *et al.* [7] were  
31  
32 281 found to cluster with *M. m. littoralis*, 16 with *M. m. lasiotis* and one with *M. m. mulatta* (Fig. 4d).  
33  
34 282 Thus, the data and results presented here provide the base date for tracing the origin of captive RMs  
35  
36 283 and the basis for the selection of appropriate animal models when testing for particular diseases, and  
37  
38 284 are thus a significant contribution to the “3Rs” principle, which aim to reduce, refine, and replace  
39  
40 285 experimental animals.  
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

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

## 286 **Conclusion**

287 We present the first description of the evolutionary history and genomic variation of geo-referenced  
288 wild RMs throughout China, including scenarios on potential functions of this variation in  
289 adaptation to local environments. This genomic resource represents a valuable contribution to the  
290 understanding of the biology and evolution of a highly successful and important biomedical research  
291 species. In particular, it is important to note that due to the difference in evolutionary history of the  
292 subspecies identified here, it can be expected that animals originating from different regions may  
293 react differently to experimental tests, and thus their background needs to be assessed beforehand  
294 [10]. Our results highlight the importance that genome typing can play in biomedical research where  
295 animal origins are uncertain, and the resources generated here provide a baseline for genomic  
296 assessment of biomedical research populations, genetic resource conservation and for refined usage  
297 of RMs in future research.

298

## 299 **Materials and Methods**

### 300 **Ethics statement**

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

### 305 **Sample Collection and Sequencing**

306 Samples from 79 individuals with information about geographic origin were collected from 17 local  
307 wildlife rescue center, which covered most of the species' range in China. Muscle samples were  
308 collected from deceased individuals and the blood samples were taken during routine physical  
309 examinations. Total genomic DNA was extracted from blood or tissue samples using standard  
310 phenol/chloroform methods. For each individual, ~3 µg DNA was sheared into fragments of 500 bp  
311 with the Covaris system. DNA fragments were then processed and sequenced using the Illumina  
312 HiSeq 2000 and 2500 platform. Raw reads were first filtered with the following criteria: (1) reads  
313 with unidentified nucleotides (N) exceeded 10% were discarded, (2) reads with the proportion of  
314 low quality base (phred quality <=5) larger than 50% were discarded. After the quality control, a  
315 total of 2,736.91 Gb of high quality sequences with 22.53 billion pair-end reads (100 or 125 bp)  
316 were generated. Furthermore, published genomic data for two individuals were download form  
317 NCBI [9,23] and filtered using the same conditions.

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

320 High-quality sequence reads were mapped to the macaque reference genome, rheMac2 [52], using  
321 the Burrows–Wheeler Aligner (BWA) (0.7.10-r789) [53]. Sequence Alignment/Map (SAM) format  
322 files were imported to SAMtools (v0.1.19) [54] for sorting and removing duplicated reads.  
323 Following mapping, we performed SNP calling using SAMTools on autosomal sites only. To obtain  
324 high-quality SNPs, we applied the calling protocol used in Chen et al [55]. The variants were filtered  
325 unless the minimum root-mean-square (RMS) mapping quality was 20. And then variants were  
326 removed if their average Phred scaled base quality was lower than 20 or the distance between the  
327 SNP was less than 5bp. Furthermore, only the variants supported by at least four reads were



1 328 presented for the subsequent analysis. Using SAMTools we discovered 55,404,179 SNPs on the  
2 329 autosomes of 81 Chinese RMs. Finally, all the SNPs were annotated by ANNOVAR (v2013-06-21)  
3  
4 330 [56] (Supplementary Table 2). For each individual the heterozygosity was calculated as  
5  
6 331 heterozygous SNP rate across the whole genome (Supplementary Table 3).  
7

### 8 332 **Genetic Diversity and Structure Analysis**

9  
10 333 A neighbor-joining (NJ) tree was constructed for the 81 individuals based on the autosomal genome  
11 334 data using the software TreeBeST. The bootstrap was set to 1,000 times to assess branch reliability,  
12 335 with the genome information of Indian RMs and *M. sylvanus* as outgroups. FigTree  
13 336 (<http://tree.bio.ed.ac.uk/software/figtree/>, v1.4.0) was used to visualize the phylogenetic tree (Fig.  
14 337 1b). Population structure analysis was performed using the software STRUCTURE 2.3.4 [27],  
15 338 which estimates individual ancestry and admixture proportions assuming  $K$  ancestral populations.  
16 339 We ran STRUCTURE five times to assess convergence and tested the number of genetic clusters  
17 340 ( $K$ ) from 2-9 (Supplementary Fig. 5). We also carried out a principle component analysis (PCA)  
18 341 using the smartPCA program from the Eigensoft package (v5.0) [57]. To determine the significance  
19 342 level of principal components, a Tracy-Widom test was done after the PCA (Supplementary Table  
20 343 4). Linkage disequilibrium for the different populations was calculated using the haploview software  
21 344 [58] with the maxdistance set as 500kb (Supplementary Fig. 13).  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37

### 38 346 **Demographic and Divergence Inference Using PSMC and Fastsimcoal2**

39 347 The PSMC model [24] was used to estimate the population histories from the individual genomes  
40 348 (sex chromosomes excluded) with the following parameters:  $-N25 -t15 -r5 -p '4+25 \times 2+4+6'$ . We  
41 349 assumed a generation length of 11 years and a mutation rate per generation ( $\mu$ ) of  $1.0 \times 10^{-8}$  [6]. To  
42 350 ensure the quality of consensus sequences, we used data of ten individuals with an average  
43 351 coverage  $>20\times$  (22.20-34.32 $\times$ ).  
44  
45  
46  
47  
48  
49

50 352 Due to the limitation of PSMC inference for recent dating, we performed the joint site  
51 353 frequency spectrum (SFS) approach implemented in *fastsimcoal2* [25] to simulate more recent  
52 354 demographic fluctuations and respective divergence times. For the five identified RM subspecies,  
53 355 eight alternative divergence scenarios describing the evolutionary relationships of these subspecies  
54 356 were tested against each other to identify the one that best supports the observed data  
55 357 (Supplementary Fig. 6). The parameters used in *fastsimcoal2* were: -N 100000 (max. number of  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 358 simulations), -L 40 (max. number of EM cycles), - M 0.001 (min. relative difference in parameter  
2 359 values for the stopping criterion). Multiple replicates with each model starting from different initial  
3  
4 360 conditions were run to ensure convergence. The best model was addressed through the maximum  
5  
6 361 value of the likelihoods and the Akaike information criterion [25]. Among all the scenarios tested,  
7  
8 362 the highest lnL and lowest AIC value was generated under the scenario 2 (Supplementary Table 5).

9  
10 363 Based on the best model (model 2) identified in the previous step, a more detailed scenario was  
11  
12 364 tested including additional population parameters of interest such as effective population sizes,  
13  
14 365 migration rates. Migration rates were ignored between subspecies which have no direct connection.  
15  
16 366 The outputs of this scenario were processed with arlsumstat to obtain distributions of various  
17  
18 367 summary statistics (Supplementary Table 6).

### 20 368 **Positive Selection**

21  
22  
23 369 To identify genomic regions that may have been subject to selection for each subspecies inhabited  
24  
25 370 in different habitats, we scanned the genome using one-to-one pair-wise comparisons between all  
26  
27 371 five subspecies. For each pairwise comparison, the differences in genetic diversity between two  
28  
29 372 subspecies were reflected by pairwise nucleotide diversity ( $\theta_\pi$ ) and the divergence in allele  
30  
31 373 frequency in two subspecies was quantified by pairwise  $F_{ST}$ . We calculated  $\theta_\pi$  for each population  
32  
33 374 and the  $F_{ST}$  between the two populations in each comparison using VCFtools [59] with a genome-  
34  
35 375 wide sliding window strategy (50-kb in length with 25-kb step). The  $F_{ST}$  values were Z-transformed  
36  
37 376 and the log value of  $\theta_\pi$  ratio ( $\theta_{\pi 2} / \theta_{\pi 1}$ ) was estimated. Putative selection targets were extracted based  
38  
39 377 on the top 5% of log-odds ratios for both Z ( $F_{ST}$ ) and log ( $\theta_\pi$ -ratio). Finally for each subspecies we  
40  
41 378 used the intersection of putative selected regions generated by all the pair-wise comparisons with  
42  
43 379 other subspecies as the candidate regions with selective pressure (i.e. consistent signatures of  
44  
45 380 selective sweeps). Genes located in these regions are expected to represent targets of selection.  
46  
47 381 Functional classification and enrichment analysis of GO categories and KEGG pathways for these  
48  
49 382 candidate genes were performed using DAVID (v6.8) [60]. The modified Fisher Exact  $P$ -value cut  
50  
51 383 off was 0.05.

### 52 384 **Genomic divergence and implication for biomedical research**

53  
54  
55 385 A total of 111 out of 55,404,179 RM SNPs analyzed in this study were successfully mapped to  
56  
57 386 human reference sequence version hg19 (GRCh37) using liftOver ([https://genome.ucsc.edu/cgi-  
58  
59 387 bin/hgLiftOver](https://genome.ucsc.edu/cgi-bin/hgLiftOver)) and were annotated as ‘disease causing’ in HGMD (version 2015.1) or pathogenic

1 388 in ClinVar (downloaded 8/01/2017) (Supplementary Table 11).

2 389

3 390 For more details of methods please see supplementary notes in Supplementary Material.

4 391

## 5 392 **Data Access**

6 393 All data generated from this study have been submitted to the NCBI Sequence Read Archive (SRA)

7 394 under BioProject PRJNA345528.

8 395

## 9 396 **Competing interests**

10 397 The authors declare that they have no competing interests.

11 398

## 12 399 **Acknowledgments**

13 400 The project was supported by the Strategic Priority Research Program of the Chinese Academy of

14 401 Sciences (Grant No. XDPB0202), NSFC (31530068 and 31471989), and National Key R&D

15 402 Program of China (2016YFC0503200). The authors thank Baoguo Li, Meng Yao, Songtao Guo,

16 403 Jiqi Lu, Zhenlong Wang, Xuelong Jiang, Tao Meng and Qihai Zhou for their help in sampling;

17 404 Daniel Pitt, Quan Kang, Qi Wu and Qi Pan for their assistance in data analysis.

18 405

## 19 406 **Author contributions**

20 407 Z. L., M. B. and M. L. conceived the study and designed the project. Z. L., X. T., P. O., X. Z. and

21 408 S. T. managed the project, performed the analyses and wrote the manuscript. Z. L., B. S. and H. X.

22 409 prepared samples. Z. L., X. T. and P. O. performed genetic analyses. Z. L., X. T., P. O., B. R., L. Z.,

23 410 G. L., Z. Y., Z. P., Z. X., C. R., M. B. and M. L. discussed the data. Z. L. and X. T. wrote the

24 411 manuscript with contributions from P. O., B. W., H. X., W. Z., C. R., M. B. and M. L.; all authors

25 412 contributed to data interpretation.

26 413

## 27 414 **Supplementary Material**

28 415 Supplementary information, figures S1-S13, tables S1-S12, and notes are available on line.

## 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. Maestriperi D. *Macchiavellian 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 Antigens*. 2011; **78**:256–62.

- 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
- 445 13. Hernandez RD, Hubisz MJ, Wheeler DA et al. Demographic histories and patterns of linkage  
446 disequilibrium in Chinese and Indian rhesus macaques. *Science* 2007; **316**:240-243.
  - 447 14. Champoux M, Higley JD, Suomi SJ. Behavioral and physiological characteristics of Indian and  
448 Chinese-Indian hybrid rhesus macaque infants. *Dev Psychobiol* 1997; **31**:49–63.
  - 449 15. Trichel AM, Rajakumar PA, Murphey-Corb M. Species-specific variation in SIV disease  
450 progression between Chinese and Indian subspecies of rhesus macaque. *J Med Primatol* 2002;  
451 **31**:171-178.
  - 452 16. Tosi AJ, Morales JC, Melnick DJ. Paternal, maternal, and biparental molecular markers provide  
453 unique windows onto the evolutionary history of macaque monkeys. *Evolution* 2003; **57**:1419-  
454 1435.
  - 455 17. Smith DG. Genetic characterization of Indian-origin and Chinese-origin rhesus macaques  
456 (*Macaca mulatta*). *Comp Med* 2005; **55**:227-230.
  - 457 18. Ferguson B, Street SL, Wright H et al. Single nucleotide polymorphisms (SNPs) distinguish  
458 Indian-origin and Chinese-origin rhesus macaques (*Macaca mulatta*). *BMC Genomics* 2007;  
459 **8**:43.
  - 460 19. Kubisch HM, Falkenstein KP, Deroche CB et al. Reproductive efficiency of captive Chinese-  
461 and Indian-origin rhesus macaque (*Macaca mulatta*) females. *Am J Primatol* 2012; **74**:174-184.
  - 462 20. Kanthaswamy S, Johnson Z, Trask JS et al. Development and validation of a SNP-based assay  
463 for inferring the genetic ancestry of rhesus macaques (*Macaca mulatta*). *Am J Primatol* 2014;  
464 **76**:1105-1113.
  - 465 21. Fooden J. Systematic review of the rhesus macaque, *Macaca mulatta* (Zimmermann, 1780).  
466 *Field Zool* 2000; **96**:1–180.
  - 467 22. Jiang X, Wang Y, Ma S. Taxonomic revision and distribution of subspecies of rhesus monkey  
468 (*Macaca mulatta*) in China. *Zool Res* 1991; **12**:241-247.
  - 469 23. Fang X, Zhang Y, Zhang R et al. Genome sequence and global sequence variation map with 5.5  
470 million SNPs in Chinese rhesus macaque. *Genome Biol* 2011; **12**:R63.
  - 471 24. Prado-Martinez J, Sudmant PH, Kidd JM et al. Great ape genetic diversity and population  
472 history. *Nature* 2013; **499**:471-475.
  - 473 25. Evanno G, Regnaut S, Goudet J. Detecting the number of clusters of individuals using the  
474 software STRUCTURE: a simulation study. *Mol Ecol* 2005; **14**:2611-2620.

- 1 475 26. Li H, Durbin R. Inference of human population history from individual whole-genome  
2 476 sequences. *Nature* 2011; **475**:493-496.
- 3  
4 477 27. Zheng B, Xu Q, Shen Y. The relationship between climate change and Quaternary glacial cycles  
5 478 on the Qinghai–Tibetan Plateau: review and speculation. *Quatern Int* 2002; **97**:93-101.
- 6  
7 479 28. Excoffier, L. Dupanloup I, Huerta-Sánchez E et al Robust demographic inference from genomic  
8 480 and SNP data. *PLoS Genet* 2013; **9**:e1003905.
- 9  
10 481 29. Owen LA, Finkel RC, Caffee MW. A note on the extent of glaciation throughout the Himalaya  
11 482 during the global Last Glacial Maximum. *Quaternary Sci Rev* 2002; **21**:147-157.
- 12  
13 483 30. Owen LA. Latest Pleistocene and Holoene glacier fluctuations in the Himalaya and Tibet.  
14 484 *Quaternary Sci Rev* 2009; **28**:2150-2164.
- 15  
16 485 31. Wu S, Luo J, Li Q et al. Ecological genetics of Chinese rhesus macaque in response to mountain  
17 486 building: all things are not equal. *PLoS ONE* 2013; **8**:e55315.
- 18  
19 487 32. Yi X, Liang Y, Huerta-Sanchez E et al. Sequencing of 50 human exomes reveals adaptation to  
20 488 high altitude. *Science* 2010; **329**:75-78.
- 21  
22 489 33. Bhatia G, Patterson N, Pasaniuc B et al. Genome-wide comparison of African-ancestry  
23 490 populations from CARE and other cohorts reveals signals of natural selection. *Am J Hum Genet*  
24 491 2011; **89**:368–381.
- 25  
26 492 34. Zhao SC, Zheng PP, Dong SS et al. Whole-genome sequencing of giant pandas provides  
27 493 insights into demographic history and local adaptation. *Nat Genet* 2013; **45**:67-71.
- 28  
29 494 35. Bergmann C. Über die Verhältnisse der Wärmeökonomie der Thiere zu ihrer Grösse. *Göttinger*  
30 495 *Studien* 1847; **3**:595–708.
- 31  
32 496 36. Zhang P, Lyu MY, Wu CF et al. Variation in body mass and morphological characters in *Macaca*  
33 497 *mulatta breviceaudus* from Hainan, China. *Am J Primatol* 2016; **78**:679-698.
- 34  
35 498 37. Zhao X, Zhang H, Lv X et al. Survey and research of morphological characters of monkeys  
36 499 (*Macaca mulatta*) in the Taihang Mountains. *J Henan Nor Uni* 1989; **62**:120-125.
- 37  
38 500 38. Kurima K, Warman ML, Krishnan S et al. A member of a family of sulfate-activating enzymes  
39 501 causes murine brachymorphism. *Proc Natl Acad Sci USA* 1998; **95**:8681-8685.
- 40  
41 502 39. Faiyaz ul Haque M, King LM, Krakow D et al. Mutations in orthologous genes in human  
42 503 spondyloepimetaphyseal dysplasia and the brachymorphic mouse. *Nat Genet* 1998; **20**:157-162.
- 43  
44 504 40. Dy P, Smits P, Silvester A et al. Synovial joint morphogenesis requires the chondrogenic action

- 1 505 of Sox5 and Sox6 in growth plate and articular cartilage. Dev Biol 2010; **341**:346-359.
- 2 506 41. Amling M, Neff L, Tanaka S et al. Bcl-2 lies downstream of parathyroid hormone-related in a
- 3 507 signaling pathway that regulates chondrocyte maturation during skeletal development. J Cell
- 4 508 Biol 1997; **136**:205-213.
- 5 509 42. Van Hengel J, Calore M, Bauce B et al. Mutations in the area composite protein  $\alpha$ T-catenin are
- 6 510 associated with arrhythmogenic rightventricular cardiomyopathy. Eur Heart J 2013; **34**:201-210.
- 7 511 43. González-Alonso J. Human thermoregulation and the cardiovascular system. Exp Physiol 2012;
- 8 512 **97**:340-346.
- 9 513 44. Ranneva SV, Pavlov KS, Gromova AV et al. Features of emotional and social behavioral
- 10 514 phenotypes of calyntenin2 knockout mice. Behav Brain Res 2017; **332**:343-354.
- 11 515 45. Horn KE, Glasgow SD, Gobert D et al. DCC expression by neurons regulates synaptic plasticity
- 12 516 in the adult brain. Cell Rep 2010; **31**:173-185.
- 13 517 46. Hori K, Hoshino M. Neuronal Migration and AUTS2 Syndrome. Brain Sci 2017; **7**:e54.
- 14 518 47. Hazen MM, Woodward AL, Hofmann I et al. Mutations of the hemophagocytic
- 15 519 lymphohistiocytosis-associated gene UNC13D in a patient with systemic juvenile idiopathic
- 16 520 arthritis. Arthritis Rheum 2008; **58**:567-570.
- 17 521 48. Procter M, Wolf B and Mao R. Forty-eight novel mutations causing biotinidase deficiency. Mol
- 18 522 Genet Metab 2016; **117**:369-372.
- 19 523 49. Munier S, Delcroix-Genete D, Carthagena L et al. Characterization of two candidate genes,
- 20 524 NCoA3 and IRF8, potentially involved in the control of HIV-1 latency. Retrovirology 2005;
- 21 525 **2**:73.
- 22 526 50. Fan ZY, Song YL. Chinese Primate Status and Primate Captive Breeding for Biomedical
- 23 527 Research in China. In: Institute for Laboratory Animal Research, National Research Council.
- 24 528 International Perspectives: The Future of Nonhuman Primate Resources. Washington DC:
- 25 529 National Academy Press. 2003.
- 26 530 51. Hao Xin. Monkey Research in China: Developing a Natural Resource. Cell 2007; **129**: 1033-
- 27 531 1036.
- 28 532 52. Gibbs RA, Rogers J, Katze MG et al. Evolutionary and biomedical insights from the rhesus
- 29 533 macaque genome. Science 2007; **316**:222-234.
- 30 534 53. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform.

1 535 Bioinformatics 2009; **25**:1754-1760.

2 536 54. Li H, Handsaker B, Wysoker A et al. The sequence alignment/map format and SAMtools.

3

4 537 Bioinformatics 2009; **25**:2078–2079.

5

6

7 538 55. Chen C, Liu Z, Pan Q et al. Genomic analyses reveal demographic history and temperate

8

9 539 adaptation of the newly discovered honey bee subspecies *Apis mellifera sinixinyuan* n. ssp.

10

11

12 540 Mol Biol Evol 2016; **33**:1337–1348.

13

14

15 541 56. Wang K, Li M, Hakonarson H. ANNOVAR: functional annotation of genetic variants from

16

17 542 high-throughput sequencing data. Nucleic Acids Res 2010; **38**:e164.

18

19 543 57. Patterson N, Price AL, Reich D. Population structure and eigenanalysis. PLoS Genet 2006;

20

21 544 **2**:e190.

22

23 545 58. Barrett JC, Fry B, Maller J et al. Haploview: analysis and visualization of LD and haplotype

24

25 546 maps. Bioinformatics 2005; **21**:263-265.

26

27 547 59. Danecek P, Auton A, Abecasis G et al. The variant call format and VCFtools. Bioinformatics

28

29 548 2011; **27**:2156-2158.

30

31 549 60. Huang da W, Sherman BT, Lempicki RA. Systematic and integrative analysis of large gene lists

32

33 550 using DAVID Bioinformatics Resources. Nat Protoc 2009; **4**:44-57.

34

35 551

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



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

552 **Figure Legends**

553 **Figure 1.** Phylogeny and population genetic structure of 81 wild Chinese RMs. (a) Geographic  
554 distribution of RMs in China (gray shadow) and the 17 sampling sites along with their subspecies  
555 assignment. (b) Neighbor-joining (NJ) tree and clustering solution inferred using STRUCTURE and  
556 displaying five populations (inferred with Evanno's  $\Delta K$  method; Supplementary Fig. 5). (c)  
557 Principal component analysis plots depicting the first two components (variance explained by PC1  
558 = 6.40% and PC2 = 5.07%).

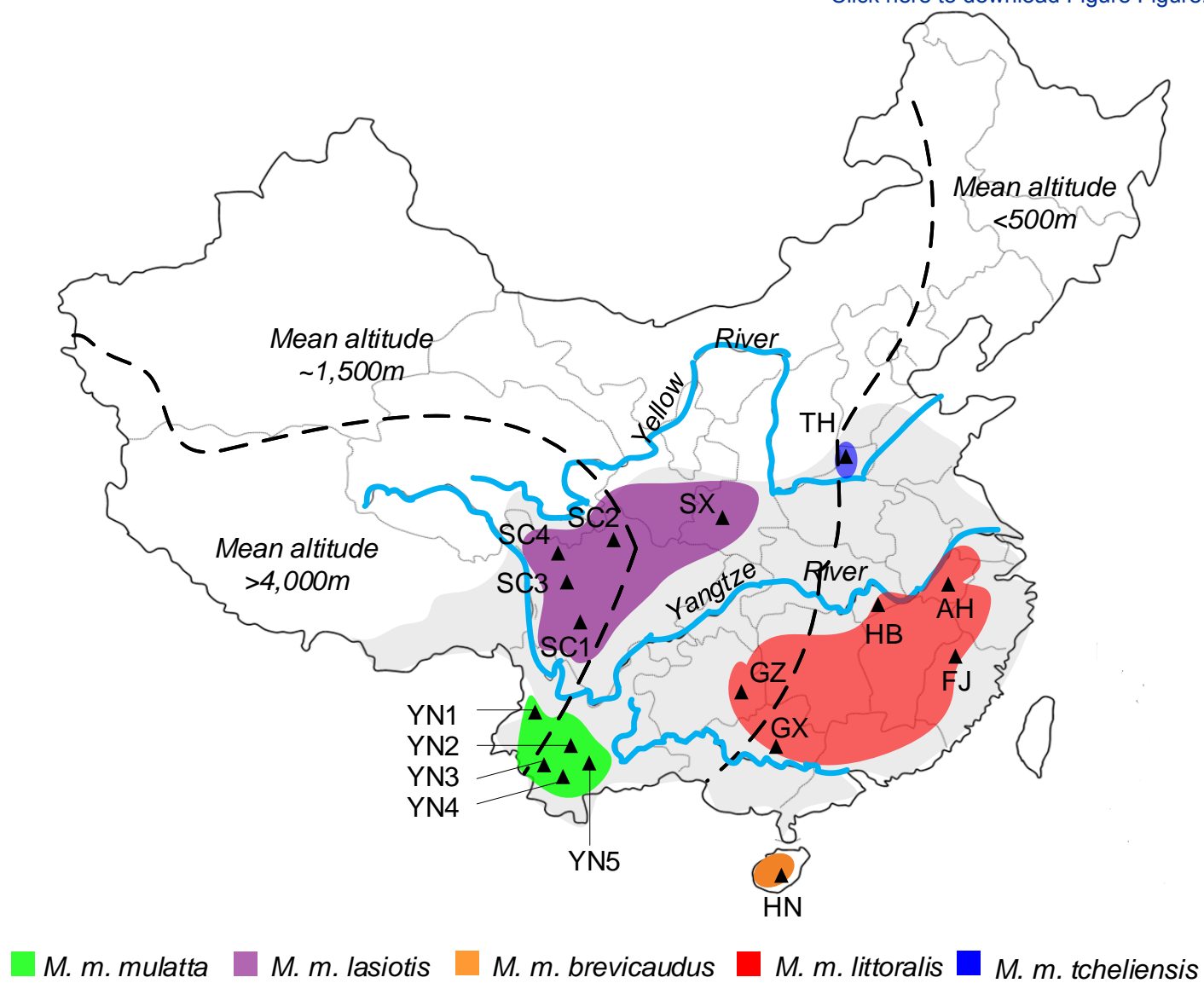
559 **Figure 2.** Demographic history and differentiation scenarios of Chinese RMs. (a) Historical changes  
560 in effective population size reconstructed using the pairwise sequential Markovian coalescent  
561 (PSMC) applied on individual whole genomes for each of the five subspecies. The generation length  
562 (g) and the neutral mutation rate per generation ( $\mu$ ) were assumed to be 11 years and  $1.08 \times 10^{-8}$ ,  
563 respectively. The Naynayxungla Glaciation (NG, 780-500 kya), Penultimate Glaciation (PG, 200-  
564 130 kya) and Last Glaciation (LG, 70-10 kya) are shaded in gray. (b) Demographic history inferred  
565 by *fastsimcoal2*. The width of the gray bars and numbers on them indicate the estimated effective  
566 population size. The arrows indicate migration patterns with the numbers above arrows indicating  
567 the average number of migrants per generation between different subspecies. Numbers at the right  
568 show the divergence times between subspecies. (c) Biogeographic scenario for RMs. Chinese RMs  
569 separates from Indian RMs ~ 162 kya [13], followed by further migration into China by the different  
570 RM subspecies indicated with arrows colored following the color key in Fig. 1a.

571 **Figure 3.** Genomic regions with selection sweep signals in RM. (a) Distribution of  $\log_2(\theta_\pi M. m.$   
572 *lasiotis*/ $\theta_\pi M. m. tcheliensis$ ) and  $Z(F_{ST})$  of 50-kb windows with 25-kb steps. Blue dots located in  
573 the selected regions requirement (corresponding to Z test  $P < 0.05$ , where  $Z(F_{ST}) \geq 1.848$  and  $\theta_\pi$   
574  $\log$ -ratio  $\geq 1.203$ ) represent selected windows for *M. m. tcheliensis*. (b) Morphological comparison  
575 between *M. m. tcheliensis* and *M. m. lasiotis*. M and F represent males and females. (c) Example of  
576 genes with selection sweep signals. *Sox5*, *Bcl2* and *Papss2* in *M. m. tcheliensis* and *Camk2g* and  
577 *Ppp3cb* in *M. m. brevicaudus*.  $F_{ST}$  and  $\theta_\pi$   $\log$ -ratio between the two subspecies are represented in red  
578 and blue, respectively. All values in figure 3c are plotted using 50-kb windows with half steps.  
579 Genome annotations are show at the bottom (black bar, coding sequences (CDS); purple and orange  
580 bar, genes). (d) SNP genotypes in putative selective sweeps containing *Sox5*, *Bcl2*, *Papss2*, *Camk2g*  
581 and *Ppp3cb*.

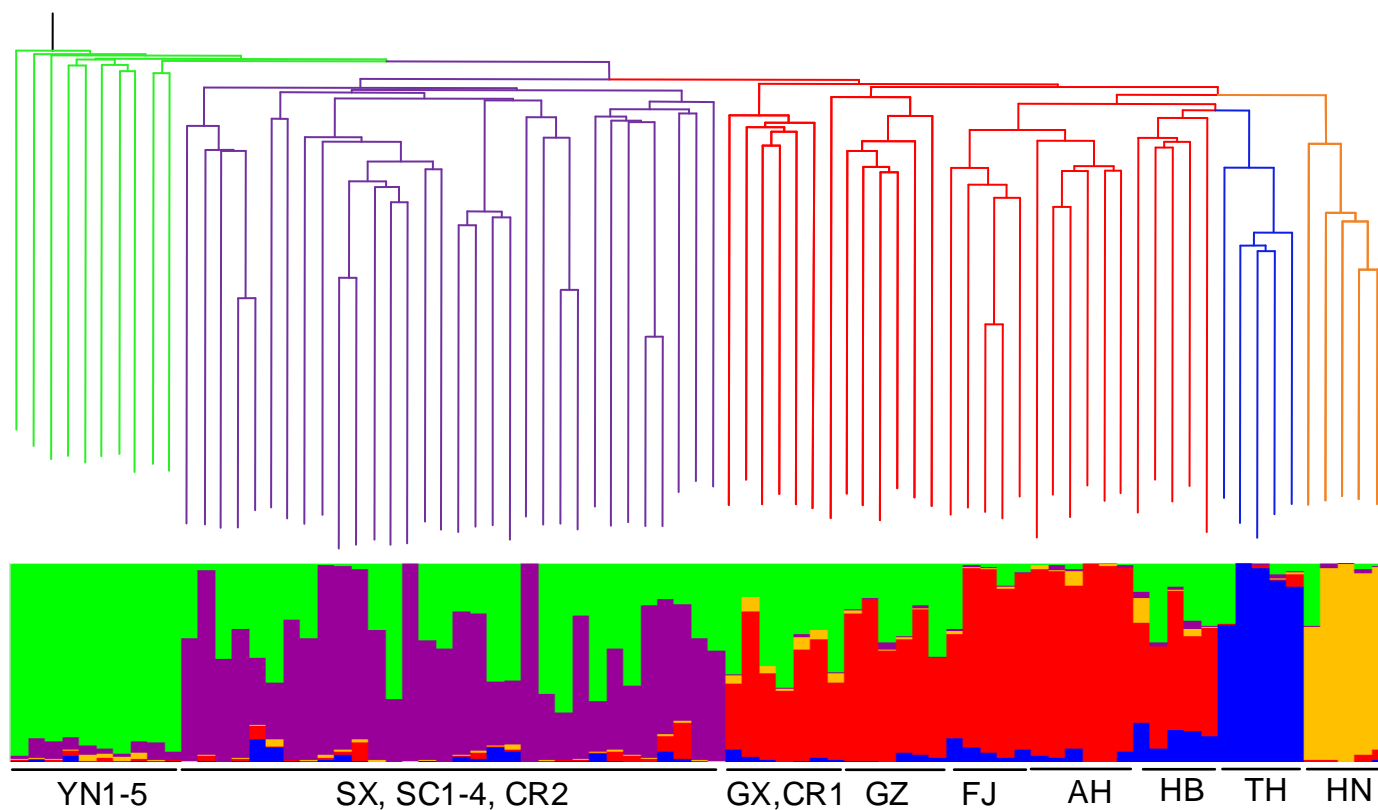
1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59

**Figure 4.** Population study of putative pathogenic SNPs in Chinese RM subspecies. (a) The site and frequency of pathogenic SNPs located in *Unc13d* and *Btd* genes. (b) Scheme of the *Ncoa3* gene in RM. The positions of nonsynonymous polymorphisms (black) and three amino-acid deletions (in red) are marked. (c) Private and shared pathogenic SNPs in Chinese RM subspecies (blue: *M. m. tcheliensis*; orange: *M. m. brevicaudus*; red: *M. m. littoralis*; green: *M. m. mulatta*; purple: *M. m. lasiotis*). The sizes of the areas are not proportional to the magnitude of the numbers. (d) NJ tree including the 81 Chinese RMs derived from this study, the 26 captive Chinese RMs from Zhong et al. [7] are indicated by blue dot.

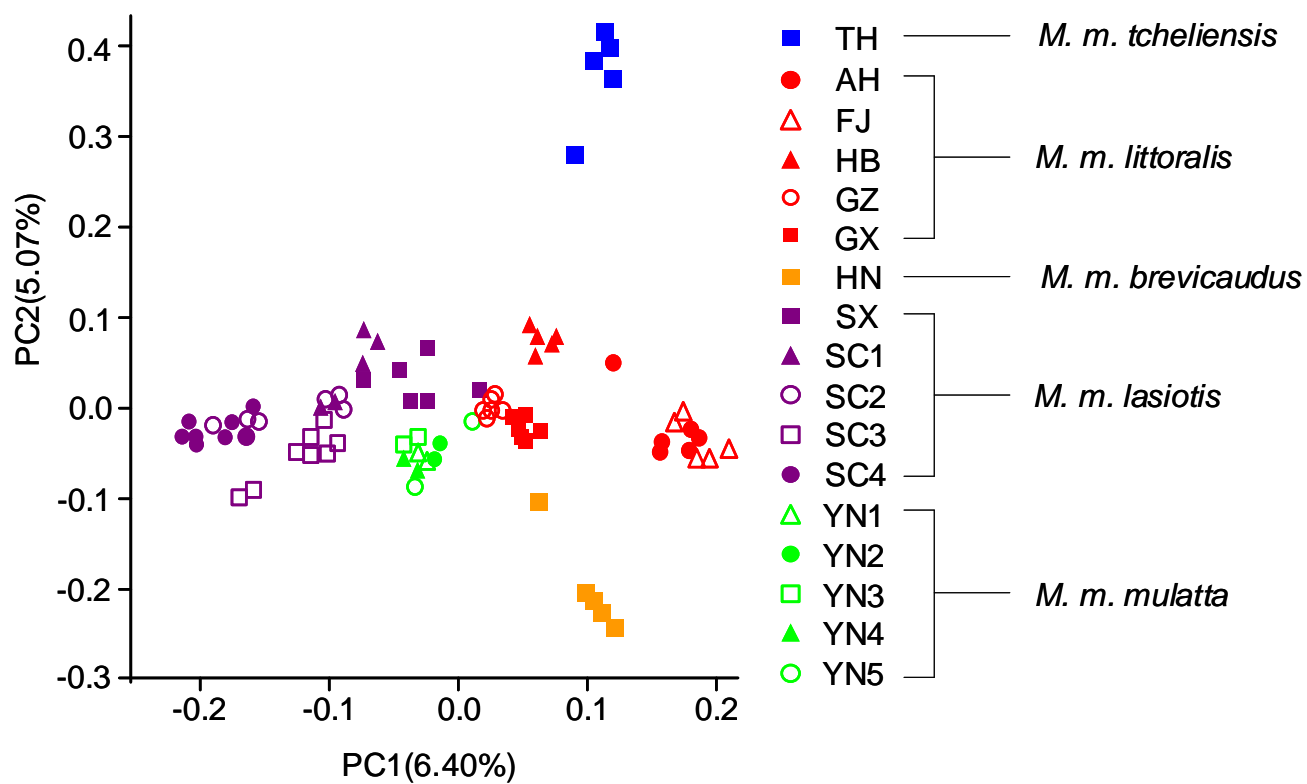
a

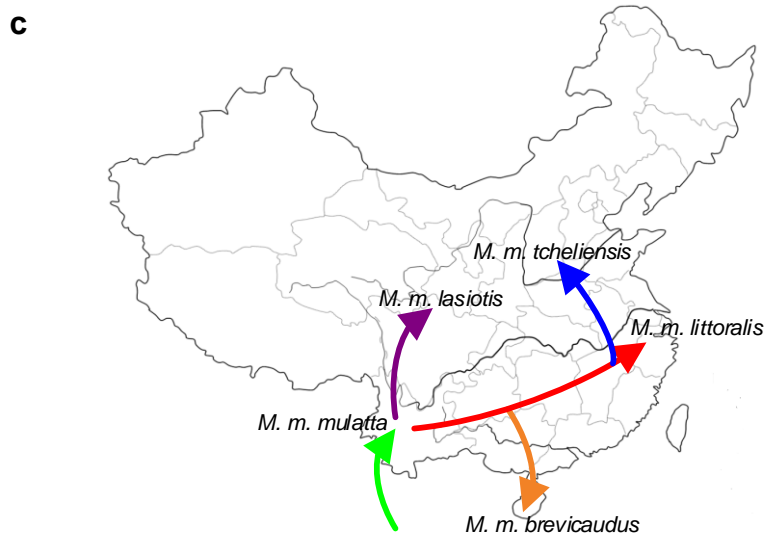
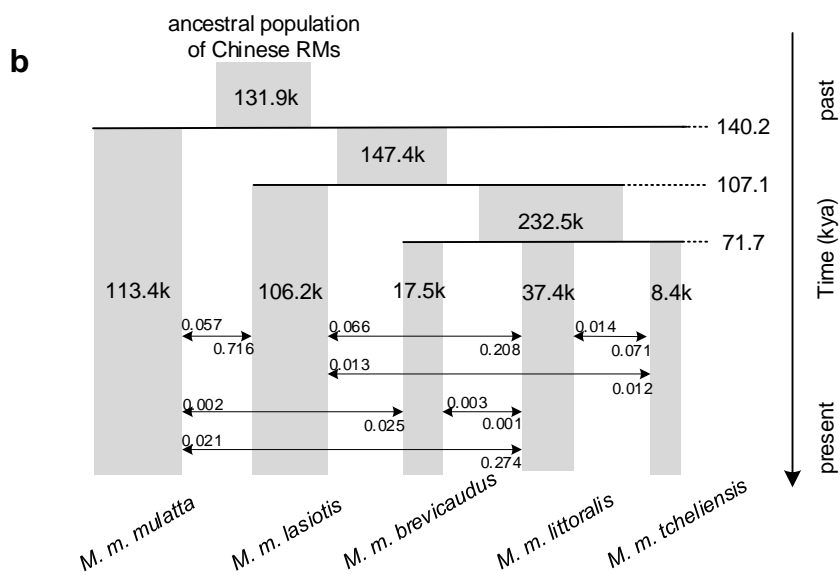
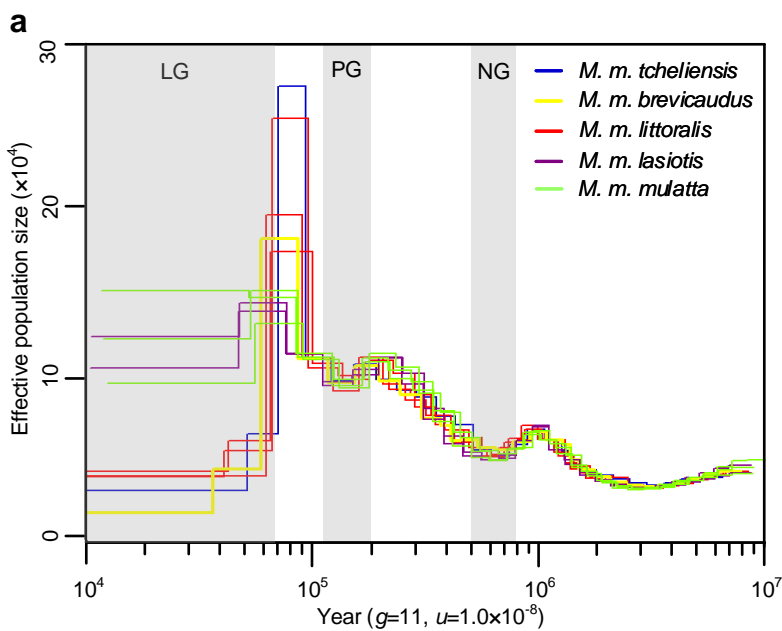


b

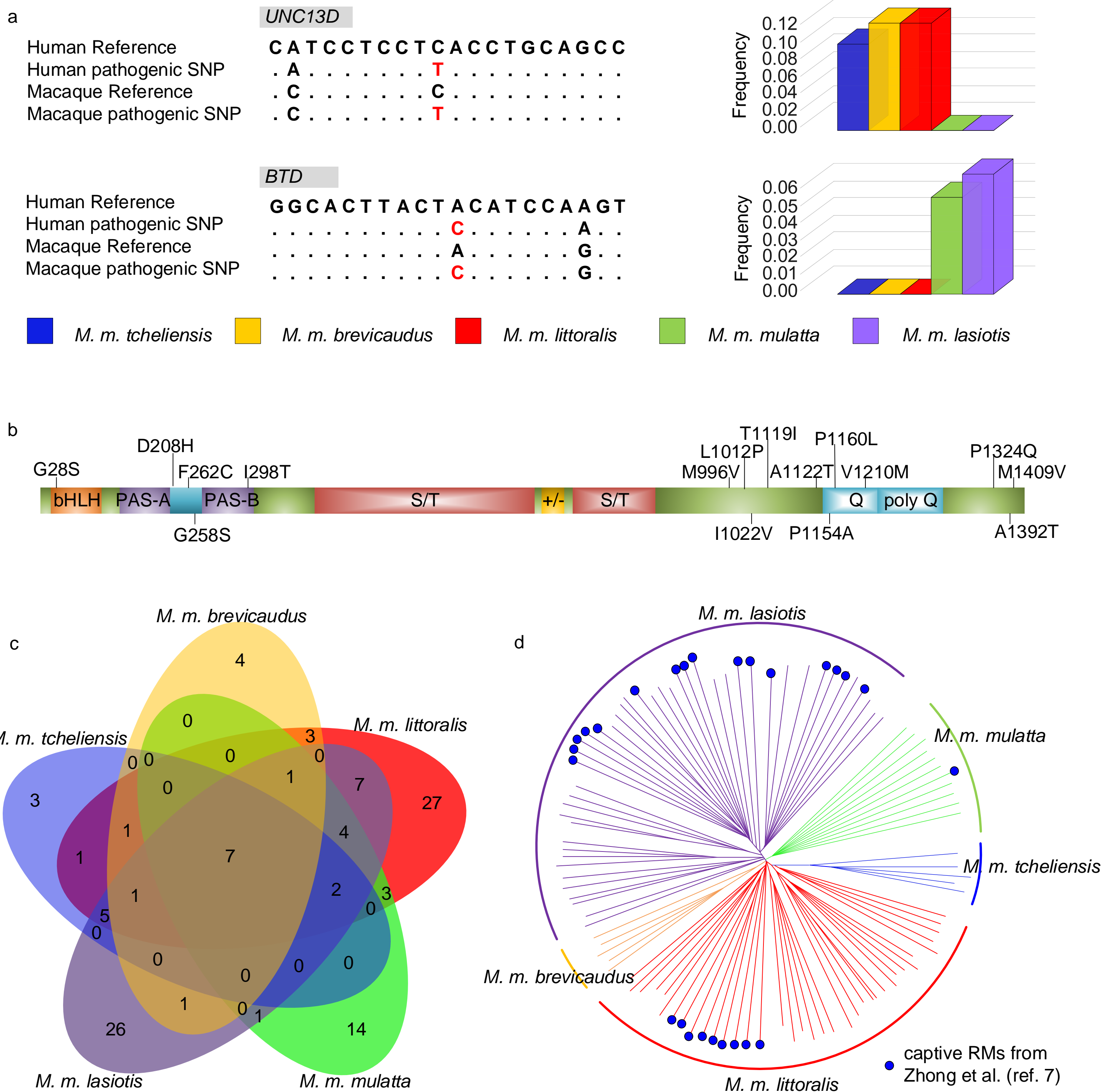


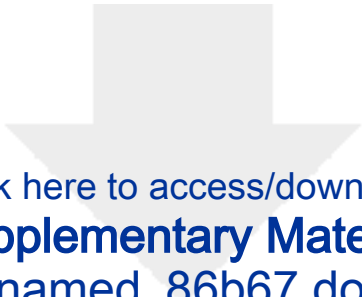
c











Click here to access/download  
**Supplementary Material**  
renamed\_86b67.docx

