

Manuscript Number:	GIGA-D-17-00160R2	
Full Title:	The genome of the Marco Polo Sheep (<i>Ovis ammon polii</i>)	
Article Type:	Data Note	
Funding Information:	National Natural Science Foundation of China (31072019)	Dr. Yutao Wang
	National Natural Science Foundation of China (31572381)	Dr. Yu Jiang
	Talents Team Construction Fund of Northwestern Polytechnical University (NWPUP)	Dr. Qiang Qiu Dr. Wen Wang
Abstract:	<p>Background: The Marco Polo Sheep (<i>Ovis ammon polii</i>), a subspecies of argali (<i>Ovis ammon</i>) which is distributed mainly in the Pamir Mountains, provides a mammalian model in which to study high-altitude adaptation mechanisms. Due to over-hunting and subsistence poaching, as well as competition with livestock and habitat loss, <i>O. ammon</i> has been categorized as an endangered species on several lists. It can have fertile offspring with sheep. Hence a high quality reference genome of the Marco Polo Sheep will be very helpful in conservation genetics and even in exploiting useful genes in sheep breeding.</p> <p>Findings: A total of 1,022.43 Gb of raw reads resulting from whole-genome sequencing of a Marco Polo Sheep were generated using an Illumina HiSeq2000 platform. The final genome assembly (2.71 Gb), which has an N50 contig size of 30.7 Kb and a scaffold N50 of 5.49 Mb. The repeat sequences identified account for 46.72% of the genome and 20,336 protein-coding genes were predicted from the masked genome. Phylogenetic analysis indicated a close relationship between Marco Polo Sheep and the domesticated sheep, and the time of their divergence was approximately 2.36 million years ago (Mya). We identified 271 expanded gene families and 168 putative positively selected genes in the Marco Polo Sheep lineage.</p> <p>Conclusions: We provide the first genome sequence and gene annotation for the Marco Polo Sheep. The availability of these resources will be of value in the future conservation of this endangered large mammal, for research into high-altitude adaptation mechanisms, for reconstructing the evolutionary history of the Caprinae and for the future conservation of the Marco Polo Sheep.</p>	
Corresponding Author:	Kun Wang, Ph. D Northwestern Polytechnical University Xi'an, Shannxi CHINA	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:	Northwestern Polytechnical University	
Corresponding Author's Secondary Institution:		
First Author:	Yongzhi Yang	
First Author Secondary Information:		
Order of Authors:	Yongzhi Yang	
	Yutao Wang	
	Yue Zhao	
	Xiuying Zhang	
	Ran Li	

	Lei Chen
	Guojie Zhang
	Yu Jiang
	Qiang Qiu
	Wen Wang
	Hongjiang Wei
	Kun Wang, Ph. D
Order of Authors Secondary Information:	
Response to Reviewers:	<p>In your next (final) version, please also remove any highlighting of changes that have been made for the purpose of peer review. Reply: We have removed all highlighting in the manuscript.</p> <p>Regarding the corresponding authors, we must insist that you list a maximum of two co-authors with this role, as explained in our instructions for authors. Reply: We have reduced one corresponding author from the manuscript. Kun wang and Hongjiang wei were now listed as the corresponding authors.</p> <p>I have only one minor comment, and that concerns the heterozygosity analysis. I think the authors are probably justified in including this section, because although it is based on just one animal, it is presumably possible for researchers studying other sheep breeds to perform similar analyses and make comparisons. I wonder whether the 14% of the genome that is highly homozygous - I think this is a better word to use than 'homogeneous' (line 152) - could be due to recent inbreeding, and that the remainder of the genome is more useful for comparing levels of diversity between Marco Polo sheep and other sheep? I'm not proposing that the authors analyse these regions in any more detail though. Reply: Thank you for the suggestion, and we agree with that recent inbreeding could lead to the highly homozygous in the genome. We have re-written related sentences in line 143-148. "The genomic regions with "low heterozygosity" state that made up 14% of the genome were highly homozygous (mean heterozygosity rate = 0.003%), which could be explained by either loss of polymorphism in endangered species [20] or recent inbreeding in some specific Macro Polo sheep individuals. More samples will be required to test whether the highly homozygous status was common in this species."</p>
Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
Experimental design and statistics	Yes
<p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	
Resources	Yes

<p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p>	
<p>Availability of data and materials</p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our Minimum Standards Reporting Checklist?</p>	<p>Yes</p>

Draft genome of the Marco Polo Sheep (*Ovis ammon polii*)

Yongzhi Yang^{1,†}, Yutao Wang^{2,3,†}, Yue Zhao^{4†}, Xiuying Zhang^{2,3}, Ran Li⁴, Lei Chen¹, Guojie Zhang⁵, Yu Jiang⁴, Qiang Qiu¹, Wen Wang¹, Hongjiang Wei^{6*}, Kun Wang^{1,*}

¹ Center for Ecological and Environmental Sciences, Northwestern Polytechnical University, Xi'an 710072, China

² College of Life and Geographic Sciences, Kashgar University, Kashgar 844000, China

³ The Key Laboratory of Ecology and Biological Resources in Yarkand Oasis at Colleges & Universities under the Department of Education of Xinjiang Uygur Autonomous Region, Kashgar University, Kashgar 844000, China

⁴ College of Animal Science and Technology, Northwest A&F University, Yangling 712100, China

⁵ Centre for Social Evolution, Department of Biology, Universitetsparken 15, University of Copenhagen, Copenhagen 2100, Denmark

⁶ Key Laboratory of Banna Miniature Inbred Pig of Yunnan Province, College of Animal Science and Technology, Yunnan Agricultural University, Kunming 650225, China

*Correspondence: wk8910@gmail.com (KW), hongjiangwei@126.com (HW)

†These authors contributed equally to this work.

1 21 **Abstract**

2
3 22 **Background:** The Marco Polo Sheep (*Ovis ammon polii*), a subspecies of argali (*Ovis*
4
5 23 *ammon*) which is distributed mainly in the Pamir Mountains, provides a mammalian
6
7 24 model to study high-altitude adaptation mechanisms. Due to over-hunting and
8
9 25 subsistence poaching, as well as competition with livestock and habitat loss, *O. ammon*
10
11 26 has been categorized as an endangered species on several lists. It can have fertile offspring
12
13 27 with sheep. Hence a high quality reference genome of the Marco Polo Sheep will be very
14
15 28 helpful in conservation genetics and even in exploiting useful genes in sheep breeding.
16
17
18
19
20
21

22 29 **Findings:** A total of 1,022.43 Gb of raw reads resulting from whole-genome
23
24 30 sequencing of a Marco Polo Sheep were generated using an Illumina HiSeq2000
25
26 31 platform. The final genome assembly (2.71 Gb), which has an N50 contig size of 30.7
27
28 32 Kb and a scaffold N50 of 5.49 Mb. The repeat sequences identified account for 46.72%
29
30 33 of the genome and 20,336 protein-coding genes were predicted from the masked
31
32 34 genome. Phylogenetic analysis indicated a close relationship between Marco Polo
33
34 35 Sheep and the domesticated sheep, and the time of their divergence was approximately
35
36 36 2.36 million years ago (Mya). We identified 271 expanded gene families and 168
37
38 37 putative positively selected genes in the Marco Polo Sheep lineage.
39
40
41
42
43
44
45
46

47 38 **Conclusions:** We provide the first genome sequence and gene annotation for the Marco
48
49 39 Polo Sheep. The availability of these resources will be of value in the future
50
51 40 conservation of this endangered large mammal, for research into high-altitude
52
53 41 adaptation mechanisms, for reconstructing the evolutionary history of the *Caprinae* and
54
55 42 for the future conservation of the Marco Polo Sheep.
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

43 **Keywords:** Marco Polo Sheep, genome assembly, annotation, evolution.

44

1 45 **Data description**

2
3 46 **Introduction to *O. ammon polii***

4
5
6 47 The Marco Polo Sheep (*Ovis ammon polii*) is a subspecies of argali (*Ovis ammon*),
7
8
9 48 named after the explorer Marco Polo. It was first described scientifically in 1841 by
10
11
12 49 Edward Blyth [1]. This subspecies is distributed mainly in the Pamir Mountains, which
13
14
15 50 consist of rugged ranges at elevations of 3,500-5,200 m [2]. The habitat of the
16
17
18 51 subspecies includes the Tajikistan Pamir Mountains [3], as well as limited regions in
19
20
21 52 China, Afghanistan, Pakistan, and Kyrgyzstan [4]. The Marco Polo Sheep species
22
23
24 53 represents a new model to study high-altitude adaptation mechanisms adopted by
25
26
27 54 mammals. Due to the sheep's impressively long horns, foreign hunters have for many
28
29
30 55 years been willing to pay large amounts of money to take part in a hunt [5] and this is
31
32
33 56 still the case today [2]. Recent studies on the status of the argali population have shown
34
35
36 57 a decline in numbers, caused mainly by over-hunting and subsistence poaching, as well
37
38
39 58 as by competition with livestock and habitat loss [6-9]. *O. ammon* has been categorized
40
41
42 59 in several protection lists, such as Appendix II of CITES (Convention on International
43
44
45 60 Trade in Endangered Species of Wild Fauna and Flora) and the IUCN (International
46
47
48 61 Union for Conservation of Nature and Natural Resources) Red List, as a vulnerable or
49
50
51 62 near threatened species. Conservation and restoration measures are therefore needed in
52
53
54 63 order to safeguard the species, and information about its genome will be a key element
55
56
57 64 in formulating an appropriate conservation strategy.

58 65
59 66 **Sequencing**

1 67 High molecular weight genomic DNA was extracted from fibroblast cells cultured from
2
3 68 the ear skin biopsy sample of a male *O. ammon polii* using a Qiagen DNA purification
4
5
6 69 kit. The sheep was originally captured from the Pamir Plateau of China and reared in
7
8
9 70 the KaShi Zoo, Kashgar Prefecture, Xinjiang Province, China. A whole-genome
10
11
12 71 shotgun sequencing strategy was applied, and a series of libraries with insert sizes
13
14
15 72 ranging from 400 base pairs (bp) to 15 kilobase pairs (kb) were constructed using the
16
17
18 73 standard protocol provided by Illumina (San Diego, CA, USA). To construct small-
19
20
21 74 insert libraries (400, 500, 600, 700 and 800 bp), DNA was sheared to the target size
22
23
24 75 range using a Covaris S2 sonicator (Covaris, Woburn, MA, USA) and ligated to
25
26
27 76 adaptors. For long-insert libraries (4, 8, 10, 12 and 15 kb), DNA was fragmented using
28
29
30 77 a Hydroshear system (Digilab, Marlborough, MA, USA). Sheared fragments were end-
31
32
33 78 labelled with biotin and fragments of the desired size were gel purified. A second round
34
35
36 79 of fragmentation was then conducted before adapter ligation. All libraries were
37
38
39 80 sequenced on an Illumina HiSeq 2000 platform (**Table S1**). A total of 1,022.43 Gb of
40
41
42 81 raw data was generated, and 624.74 Gb of clean data was retrieved after removal of
43
44
45 82 duplicates, contaminated reads (reads with adaptor sequence) and low quality reads
46
47
48 83 using the sickle software tool (<https://github.com/najoshi/sickle>) with a quality
49
50
51 84 threshold of 10 and a length threshold of 50. We further corrected the short-insert library
52
53
54 85 reads using SOAPec [10], a k-mer-based error correction package.

55 87 **Evaluation of genome size**

56
57
58 88 Approximately 65 Gb clean reads were randomly selected from all short libraries to
59
60
61
62
63
64
65

1 89 estimate the genome size using the k-mer-based method and the formula: $G = k$ -
2
3
4 90 mer_number/k-mer_depth. In this study, a total of 52,413,427,492 k-mers were
5
6 91 generated and the peak k-mer depth was 17. The genome size was estimated to be
7
8
9 92 approximately 3 Gb (**Table S2** and **Fig. S1**) and all the clean data correspond to a
10
11
12 93 coverage of ~ 208-fold.

13
14
15 94

16 17 95 ***De novo* genome assembly**

18
19
20 96 The assembly was performed using Platanus v1.2.4 (Platanus, RRID:SCR_015531)
21
22
23 97 [11], which is well suited to high-throughput short reads and heterozygous diploid
24
25
26 98 genomes. Briefly, error-corrected paired-end reads (insert size < 2 kb) were input for
27
28
29 99 contig assembly with the default parameters. Next, all cleaned paired-end (insert size <
30
31
32 100 2 kb) reads and mate-paired (insert size > 2 kb) reads were mapped onto the contigs for
33
34
35 101 scaffold building, using default parameters except that the minimum number of links (-
36
37
38 102 l) was set to 10 in order to minimize the number of scaffolding errors. After gap filling
39
40
41
42 103 by Platanus, the gaps that still remained in the resulting scaffolds were closed using
43
44
45 104 GapCloser (GapCloser, RRID:SCR_015026) [10]. The final *de novo* assembly for the
46
47
48 105 Marco Polo Sheep has a total length of 2.71 Gb, including 116.91 Mb (4.3%) unknown
49
50
51 106 bases. The assembly is slightly larger than that of the domestic sheep (*Ovis aries*,
52
53
54 107 Oar_v3.1, 2.61 Gb) [12] and smaller than that of the domestic goat (*Capra hircus*,
55
56
57 108 ARS1, 2.92 Gb) [13]. The N50s for contigs and scaffolds of the Marco Polo Sheep
58
59
60 109 genome are, respectively, 30.8 kb and 5.5 Mb (**Table S3**). The assembled scaffolds
61
62
63 110 represented ~ 88% of the estimated genome size, and the GC content was 41.9%,
64
65

111 similar to those of sheep (41.9%) and goat (41.5%) (**Fig. S2**).

112 We assessed the quality of the genome assembly with respect to base-level accuracy,
113 integrity, and continuity. More than 99.65% of the short insert paired-end reads could
114 be mapped to the assembly and more than 98.35% of the sequence have a coverage
115 depth greater than 20-fold (**Table S4**), thus the assembly is of high level of single-base
116 accuracy. A core eukaryotic genes (CEG) mapping approach (CEGMA,
117 RRID:SCR_015055, v2.5 [14]) dataset comprising 248 CEGs was used to evaluate the
118 completeness of the draft: 93.55% (232/248) of genes were completely or partially
119 covered in the assembled genome (**Table S5**). Alongside this, we also used the BUSCO
120 v2.0.1 (BUSCO, RRID:SCR_015008 [15]; the representative mammal gene set
121 mammalia_odb9, which contains 4,104 single-copy genes that are highly conserved in
122 mammals) software package to assess the quality of the genome assembly generated.
123 The resulting BUSCO value was 95.9%, containing C: 92.5% [S: 91.3%, D: 1.2%], F:
124 3.4%, M: 4.1%, n: 4104 (C: complete [D: duplicated], F: fragmented, M: missed, n:
125 genes) (**Table S6**). Both the CEGMA and the BUSCO scores are comparable to those
126 for sheep (Oar_v3.1) and domestic goat (ARS1 and CHIR_1.0), which are known for
127 their high quality as the references genomes of two important livestock animals,
128 suggesting our Marco Polo Sheep assembly is of high quality and quite complete.
129 Finally, to evaluate the trade-off between the contiguity and correctness of our assembly,
130 we applied the feature-response curve (FRC) method [16], which predicts the
131 correctness of an assembly by identifying ‘features’ representing potential errors or
132 complications on each *de novo* assembled scaffold during the assembly process. The

1 133 FRC curve was calculated for the Marco Polo Sheep, sheep, taurine cattle and two
2
3
4 134 versions of goat assemblies (**Fig. S3**). We found that the curve for our assembly was
5
6 135 similar to that for the sheep and the two goat assemblies, with taurine cattle slightly
7
8
9 136 different from the others, indicating the level of contiguity and correctness of the Marco
10
11
12 137 Polo Sheep genome assembly is comparable to those of sheep and goat.

13
14 138 We mapped the reads from short-insert length libraries to the Marco Polo Sheep
15
16
17 139 reference genome with BWA (BWA, RRID:SCR_010910) [17] and performed variant
18
19
20 140 calling with SAMtools v0.1.19 (SAMTOOLS, RRID:SCR_002105) [18]. Applying
21
22
23 141 strict quality control and filtering, we obtained a total of 3.5 million SNVs (**Table S7**)
24
25
26 142 and noted that the heterozygosity rate (0.14%) was lower than that estimated for sheep
27
28
29 143 (Oar_v3.1, 0.2%) and similar with that of goat (CHIR_1.0, 0.13%) [12]. We further
30
31
32 144 assessed the distribution of heterozygosity ratio of non-overlapping 50K windows (Fig.
33
34 145 S4). We assume that the heterozygosity on the genome can be divided into three states
35
36
37 146 (low/normal/high) and applied Hidden Markov model with depmixS4 package [19] in
38
39
40 147 R to infer the state of each window. The genomic regions with “low heterozygosity”
41
42
43 148 state that made up 14% of the genome were highly homozygous (mean heterozygosity
44
45
46 149 rate = 0.003%), which could be explained by either loss of polymorphism in endangered
47
48
49
50
51 150 species [20] or recent inbreeding in some specific Macro Polo sheep individuals. More
52
53
54 151 samples will be required to test whether the highly homozygous status was common in
55
56
57 152 this species. 156 genes were overlap of more than half length with the low
58
59
60 153 heterozygosity regions and the GO enrichment analysis shows that no GO category was
61
62
63 154 significant enriched (Table S8). A total of 384,018 insertions and deletions (InDels)

1 155 (**Table S9**) were obtained. Similar to the findings of previous studies on yak [21] and
2
3
4 156 wisent [22], the InDels in the coding regions were enriched for sizes that are multiples
5
6 157 of three bases (**Fig. S5**).
7
8
9 158

11 159 **Annotation**

10
11
12
13
14 160 The transposable elements present in Marco Polo sheep sequences were identified using
15
16
17 161 a combination of *de novo* and homology-based approaches. Transposable elements
18
19
20 162 were identified at both the DNA and the protein levels, based on known sequences
21
22
23 163 contained within the DNA repeat database (RepBase v21.01) [23], using RepeatMasker
24
25 164 v4.0.5 (RepeatMasker, RRID:SCR_012954) [24] and RepeatProteinMask
26
27
28 165 (v4.0.5, a package within RepeatMasker). For the *de novo* prediction, firstly
29
30
31 166 RepeatModeler V1.0.8 (RepeatModeler, RRID:SCR_015027) was employed to
32
33
34 167 construct a *de novo* repeat library, then RepeatMasker was used to identify repeats using
35
36
37 168 both the *de novo* repeat database and RepBase. We then combined the *de novo*
38
39
40 169 prediction and the homolog prediction of transposable elements according to the
41
42
43 170 coordination in the genome. Tandem repeats were annotated with RepeatMasker and
44
45 171 Tandem Repeats Finder (TRF, V4.07) [25]. In summary, a total of 0.87% tandem
46
47
48 172 repeats and 46.60% transposable elements were identified in the Marco Polo sheep
49
50
51 173 assembly, with LINEs constituting the greatest proportion, 72.48% of all repeats, and
52
53 174 SINEs making up 24.09% of all repeats (**Table S10** and **Table S11**).

54
55
56 175 We used homology-based and *de novo* prediction to annotate protein coding genes.

57
58 176 For homology-based prediction, protein sequences from 5 different species (*Bos taurus*,

1 177 *Equus caballus*, *Homo sapiens*, *Ovis aries*, *Sus scrofa*) (**Table S12**) were mapped onto
2
3
4 178 the repeat-masked Marco Polo sheep genome using TblastN with an E-value cutoff of
5
6 179 1e-5; the aligned sequences as well as the corresponding query proteins were then
7
8
9 180 filtered and passed to GeneWise (GeneWise, RRID:SCR_015054) [26] to search for
10
11 181 accurately spliced alignments. For *de novo* prediction, we first randomly selected 1500
12
13
14 182 full-length genes from the results of homology-based prediction to train the model
15
16
17 183 parameters for Augustusv3.2.1 (Augustus: Gene Prediction, RRID:SCR_008417) [27]
18
19
20 184 and geneid v1.4.4 [28]. GenScan [29], Augustus v3.2.1 [27] and geneid v1.4.4 [28]
21
22
23 185 were then used to predict genes based on the training set of human and Marco Polo
24
25 186 Sheep genes. We used EvidenceModeler software (EVM, version 1.1.1) to integrate
26
27
28 187 the genes predicted by the homology and *de novo* approaches and generated a consensus
29
30
31 188 gene set (Table S13). The final gene set was produced by removing low-quality genes
32
33
34 189 of short length (proteins with fewer than 50 amino acids) and/or exhibiting premature
35
36
37 190 termination. The final total gene set consisted of 20,336 genes, and the number of genes,
38
39
40 191 gene length distribution and exon number per gene were similar to those of other
41
42
43 192 mammals, while the intron length was slightly larger than goat (CHIR_1.0), sheep
44
45 193 (Oar_v3.1) and taurine cattle (UMD3.1) (Table S14 and Fig. S6, S7). The repeat content
46
47
48 194 was annotated by RepeatMasker v4.0.5 [24] with unified parameters for Macro Polo
49
50
51 195 sheep, domestic sheep and goat. We found that there were more LINE sequences in the
52
53
54 196 intron regions of Marco Polo sheep than the other species, suggesting that transposon
55
56 197 insertions might have contributed to intron length increasing (Fig. S8). 92.55% of all
57
58
59 198 the predicted genes could be annotated using five protein databases: InterPro
60
61
62
63
64
65

199 (InterPro, RRID:SCR_006695) (87.17%), GO (Gene ontology, 70.99%), Swiss-Prot
200 (91.67%), TrEMBL (92.33%) and KEGG (KEGG, RRID:SCR_012773) (Kyoto
201 Encyclopedia of Genes and Genomes, 57.25%) (**Table S15**). In addition, we identified
202 2,978 noncoding RNAs in the Marco Polo Sheep genome (**Table S16**).

203

204 **Genome evolution**

205 Firstly, large-scale variations among Marco Polo Sheep, sheep and goat were identified
206 by the synteny analysis using the program LAST (LAST, RRID:SCR_006119) [30]. A
207 total of 2.29/2.30/2.40 Gb 1:1 alignment sequences were generated for, respectively
208 Marco Polo Sheep vs sheep (Oar_v3.1), Marco Polo Sheep vs goat (ASR1), sheep vs
209 goat, covering more than 88.55% of each genome (**Table S17** and **Fig. S9**). The
210 sequences present on sheep/goat autosomes were well covered (average values:
211 89.65%/89.88%) by the synteny alignment, whereas only 66.09%/63.03% were
212 covered in the case of chromosome X. The scaffolds of the Marco Polo Sheep genome
213 that aligned to the sex chromosomes were also more fragmented. The divergence
214 between Marco Polo Sheep vs sheep (Oar_v3.1), Marco Polo Sheep vs goat (ASR1),
215 sheep vs goat was 0.7%, 2.2%, 2.3%, respectively, corresponding to their relatedness
216 (**Table S17** and **Fig. S10**). Although Marco Polo Sheep, sheep and goat showed good
217 synteny alignments, there are large numbers of inter-chromosomal rearrangements
218 between pairs of them (**Fig. S11** and **S12**). By comparing Marco Polo Sheep and
219 sheep/goat genomes we identified 11,756/6,026 inter-chromosomal, intra-
220 chromosomal, or inversion breakpoints (edges of transposition events) (**Table S18**),

1 221 which may have been caused by the real translocations events between them as they
2
3 222 have a different karyotype, errors in the assembly of the genomes or erroneous synteny
4
5
6 223 alignments (false positives and false negatives). However, at this stage it is difficult to
7
8
9 224 distinguish between possible artifactual and real effects. The breakpoint distributions
10
11
12 225 were significantly enriched in repeat regions (**Fig S13a**), which are susceptible to
13
14
15 226 rearrangements but also to assembly or alignment errors. Longer scaffolds were found
16
17
18 227 to harbor fewer breakpoints (**Fig. S13b**). Single molecule sequencing with unbiased
19
20
21 228 long reads will be the best way of identifying large-scale variation.

22
23 229 To analyze gene families, we downloaded the protein sequences of eight additional
24
25
26 230 species (Opossum, human, dog, horse, pig, taurine cattle, goat and sheep) from Ensembl
27
28
29 231 (Ensembl, RRID:SCR_002344) [31] and GigaDB (GigaDB, RRID:SCR_004002) [32]
30
31
32 232 (**Table S12**). The consensus gene set for the above eight species and Marco Polo Sheep
33
34
35 233 were filtered to retain the longest CDS (coding sequence) for each gene, removing CDS
36
37
38 234 with premature stop codons and those protein sequences < 50 amino acids in length,
39
40
41 235 resulting in a dataset of 188,359 protein sequences, which was used as the input file for
42
43
44 236 OrthoMCL (OrthoMCL DB: Ortholog Groups of Protein Sequences,
45
46
47 237 RRID:SCR_007839) [33]. A total of 17,578 OrthoMCL families were built utilizing an
48
49
50 238 effective database size of all-to-all BLASTP strategy with an E-value of 1e-5 and a
51
52
53 239 Markov Chain Clustering default inflation parameter (**Table S19** and **Fig. 1a**). We
54
55
56 240 identified 155 gene families that were specific to the Marco Polo Sheep when
57
58
59 241 comparing with taurine cattle, sheep, goat and horse (**Fig. 1b**), and detected 271 gene
60
61
62 242 families that have expanded in the Marco Polo Sheep lineage using CAFÉ

1 243 (Computational Analysis of gene Family Evolution, v4.0.1) [34] (**Fig. 1a**). The
2
3 244 expanded gene families were enriched in 38 GO categories and their functions were
4
5
6 245 mainly associated with response to stimulus, cell adhesion, G-protein coupled receptor
7
8
9 246 and enzyme activity (**Table S20**).

10
11 247 Next, we selected 5,788 single-copy gene families from the above-mentioned 9
12
13 248 mammalian species and used PRANK v3.8.31 [35] with the codon option to align the
14
15
16 249 CDS from each single-copy gene family. 4D-sites (fourfold degenerate sites) were
17
18
19 250 extracted from all the single-copy genes and used to construct a phylogenetic tree with
20
21
22 251 the GTR+G+I model in RAxML v7.2.8 (RAxML, RRID:SCR_006086) [36] (**Fig. S14**).
23
24
25 252 The divergence time of each node was estimated by the PAML (PAML,
26
27
28 253 RRID:SCR_014932) MCMCtree program v4.5 [37] and calibrated against the timing
29
30
31 254 of the divergence of the opossum and human (124.6-134.8 Mya), human and taurine
32
33
34 255 cattle (95.3-113 Mya), taurine cattle and pig (48.3-53.5 Mya), and taurine cattle and
35
36
37 256 goat (18.3-28.5 Mya) [38]. The convergence was checked by Tracer v1.5 [39] and
38
39
40 257 confirmed by two independent runs. The phylogenetic analysis showed that the Marco
41
42
43 258 Polo Sheep has a closer relationship with sheep than with other mammals and that the
44
45
46 259 divergence time between them is about 2.36 (1.94-2.61) Mya (**Fig. 1a**).

47
48 260 We further used the free ratio model to calculate the average Ka/Ks values and the
49
50
51 261 branch-site likelihood ratio test to identify positively selected genes in the Marco Polo
52
53
54 262 Sheep lineage. A total of 10,353 high confidence single-copy genes were identified by
55
56
57 263 InParanoid and MultiParanoid within the human, dog, taurine cattle, goat, sheep and
58
59
60 264 Marco Polo Sheep. We found that the Marco Polo Sheep has a regular level of the
61
62
63
64
65

1 265 average Ka/Ks values, but containing more outliers (**Fig. 1c**). A total of 168 positively
2
3 266 selected genes were identified in the Marco Polo Sheep lineage (Table S21), and six of
4
5
6 267 them were orthologous with high altitude adaptation related genes (*IDE*, *IGF1*, *P2RX3*,
7
8
9 268 *PHF6*, *PROX1* and *RYR1*) identified in Tibet wild boar [40]. Two genes were
10
11 269 associated with hypoxia response: the ryanodine receptor protein encoded by *RYR1*
12
13 270 (*Ryanodine Receptor 1*) was located in the pulmonary artery smooth muscle cells,
14
15 271 which could subserve coupled O₂ sensor and NO regulatory functions to response to
16
17 272 the tissue hypoxic decrease [41]; *P2RX3* (*Purinergic Receptor P2X, Ligand-Gated Ion*
18
19 273 *Channel, 3*), is reported as a potential new target for the control of human hypertension,
20
21 274 which could reduce the arterial pressure and basal sympathetic activity and normalize
22
23 275 carotid body hyperreflexia in conscious rats with hypertension during *P2RX3*
24
25 276 antagonism [42]. Four genes were related with energetic metabolism: *IGF1* (*Insulin-*
26
27 277 *like Growth Factor 1*) encodes the growth-promoting polypeptide mainly involved in
28
29 278 the body growth and differentiation and as well as the glucose, lipid and protein
30
31 279 metabolism [43]; *IDE* (*Insulin Degrading Enzyme*) encodes a zinc metallopeptidase
32
33 280 that degrades intracellular insulin, which could accelerates glycolysis, pentose
34
35 281 phosphate cycle, and glycogen synthesis in liver [44]; *PHF6* (*PHD Finger Protein 6*)
36
37 282 encodes a protein with two PHD-type zinc finger domains and its function was
38
39 283 associated with Börjeson-Forssman-Lehmann syndrome, which is one of the syndromic
40
41 284 obesities in humans [45]; the protein encoded by *PROX1* (*Prospero Homeobox 1*) could
42
43 285 occupy promoters of metabolic genes on a genome-wide scale to control of energy
44
45 286 homeostasis [46]. In addition, the other identify PSGs may also be associated to high
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 287 altitude adaptation, while there are rare literature data on the function of them. Further
2
3 288 studies will be required to clarify the roles of these genes in high altitude tolerance.
4
5

6 289 Finally, we inferred the demographic history of the Marco Polo Sheep using the
7
8 290 Pairwise Sequentially Markovian Coalescent (PSMC) model [47]. Consensus
9
10 291 sequences were obtained using SAMtools v0.1.19 [18] and divided into non-
11
12 292 overlapping 100 bp bins. The analysis was performed with the following parameters: -
13
14 293 N25 -t15 -r5 -p '4+25×2+4+6'. PSMC modeling was done using a bootstrapping
15
16 294 approach, with sampling performed 100 times to estimate the variance of the simulated
17
18 295 results. The effective population size (N_e) of Marco Polo Sheep shows a peak at ~1 Mya
19
20 296 followed by two distinct declines. The most recent decline involved at least a sevenfold
21
22 297 decrease in N_e , and occurred ~ 60,000 years ago (Fig. S15).
23
24
25
26
27
28
29
30

31 298

32 33 34 299 **Conclusion**

35
36 300 In summary, the novel genome data generated in this work will provide a valuable
37
38 301 resource for studying high-altitude adaptation mechanisms within mammals and for
39
40 302 investigating the evolutionary histories of the *Caprinae*, and it will have relevance for
41
42 303 the future conservation of the Marco Polo Sheep.
43
44
45
46

47 304

48 49 50 305 **Availability of supporting data**

51
52 306 The sequencing reads of each sequencing library have been deposited at NCBI with the
53
54 307 Project ID: PRJNA391748, Sample ID: SAMN07274464, and the Genome Sequence
55
56 308 Archive [48] in BIG Data Center [49], Beijing Institute Genomics (BIG), Chinese
57
58
59
60
61
62
63
64
65

1 309 Academy of Science, under accession number PRJCA000449 (publicly accessible at
2
3 310 <http://bigd.big.ac.cn/gsa>). The assembly and annotation of the Marco Polo Sheep
4
5
6 311 genome are available in the *GigaScience* database GigaDB (GigaDB,
7
8 312 RRID:SCR_004002) [50]. Supplementary figures and tables are provided in Additional
9
10
11 313 file 1.

12
13
14 314

15 16 17 315 **Competing interests**

18
19
20 316 The authors declare that they have no competing interests.

21
22
23 317

24 25 318 **Authors' contributions**

26
27
28 319 KW and WW conceptualized the research project. KW, WW and HW designed analytic
29
30 320 strategy and coordinated the project. YW, HW and WW collected the samples and led
31
32 321 the genome sequencing. YY and KW led the bioinformatics analysis. YY, YW and YZ
33
34 322 generated the genome assembly and the genome annotation. YY, RL and LC finished
35
36 323 the synteny analysis. YW, YZ and GZ performed the gene family construction and the
37
38 324 phylogeny analysis. YY and QQ detected the PSGs and carried out data submission.
39
40 325 YY, WW and KW wrote the paper. All authors read and approved the final manuscript.

41
42
43
44
45
46
47 326

48 49 50 327 **Acknowledgements**

51
52
53 328 This study was supported by research grants from the National Natural Science
54
55 329 Foundation of China (No. 31072019 and No. 31572381), and Talents Team
56
57
58 330 Construction Fund of Northwestern Polytechnical University (NWPU) to QQ and WW.

1 331 We thank Nowbio Biotech Inc., Kunming, China for the remarkable work on DNA
2
3 332 libraries constructions and the assistance during the genome sequencing.
4
5
6 333
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

References

1. Dohner JV. The encyclopedia of historic and endangered livestock and poultry breeds. Yale University Press. 2001:p. 514.
2. Schaller GB and Kang A. Status of Marco Polo sheep *Ovis ammon polii* in China and adjacent countries: conservation of a Vulnerable subspecies. *Oryx*. 2008;42 1:100-6. doi:10.1017/S0030605308000811.
3. Breu TMH, Hans The Tajik Pamirs: Challenges of sustainable development in an isolated mountain region. Centre for Development and Environment (CDE), University of Berne: Berne, Switzerland. 2003:p. 80.
4. Valdez R, Michel S, Subbotin A and Klich D. Status and population structure of a hunted population of Marco Polo Argali *Ovis ammon polii* (Cetartiodactyla, Bovidae) in Southeastern Tajikistan. *Mammalia*. 2016;80 1:49-57. doi:10.1515/mammalia-2014-0116.
5. Harris RB. Ecotourism versus trophy hunting: incentives toward conservation in Yeniugou, Tibetan Plateau, China. *Integrating People and Wildlife for a Sustainable Future* (eds JA Bissonette & PR Krausman). 1995:228-34.
6. Harris RB and Reading R. *Ovis ammon*. The IUCN Red List of Threatened Species 2008: e.T15733A5074694. <http://dx.doi.org/10.2305/IUCN.UK.2008.RLTS.T15733A5074694.en>. Downloaded on 02 May 2017. 2008.
7. Shackleton DM. Wild sheep and goats and their relatives. 1997.
8. Nowak R. Court upholds controls on imports of argali trophies. *Endangered Species Technical Bulletin*. 1993;18 4:11-2.
9. Shrestha R and Wegge P. Wild sheep and livestock in Nepal Trans-Himalaya: coexistence or competition? *Environmental Conservation*. 2008;35 02:125-36.
10. Luo R, Liu B, Xie Y, Li Z, Huang W, Yuan J, et al. SOAPdenovo2: an empirically improved memory-efficient short-read de novo assembler. *Gigascience*. 2012;1 1:18. doi:10.1186/2047-217X-1-18.
11. Kajitani R, Toshimoto K, Noguchi H, Toyoda A, Ogura Y, Okuno M, et al. Efficient de novo assembly of highly heterozygous genomes from whole-genome shotgun short reads. *Genome Res*. 2014;24 8:1384-95. doi:10.1101/gr.170720.113.
12. Jiang Y, Xie M, Chen WB, Talbot R, Maddox JF, Faraut T, et al. The sheep genome illuminates biology of the rumen and lipid metabolism. *Science*. 2014;344 6188:1168-73. doi:10.1126/science.1252806.
13. Bickhart DM, Rosen BD, Koren S, Sayre BL, Hastie AR, Chan S, et al. Single-molecule sequencing and chromatin conformation capture enable de novo reference assembly of the domestic goat genome. *Nat Genet*. 2017;49 4:643-50. doi:10.1038/ng.3802.
14. Parra G, Bradnam K, Ning Z, Keane T and Korf I. Assessing the gene space in draft genomes. *Nucleic Acids Res*. 2009;37 1:289-97. doi:10.1093/nar/gkn916.
15. Simao FA, Waterhouse RM, Ioannidis P, Kriventseva EV and Zdobnov EM. BUSCO: assessing genome assembly and annotation completeness with

- 377 single-copy orthologs. *Bioinformatics*. 2015;31 19:3210-2.
 378 doi:10.1093/bioinformatics/btv351.
- 379 16. Vezzi F, Narzisi G and Mishra B. Reevaluating assembly evaluations with
 380 feature response curves: GAGE and assemblathons. *PLoS One*. 2012;7
 381 12:e52210. doi:10.1371/journal.pone.0052210.
- 382 17. Li H. Aligning sequence reads, clone sequences and assembly contigs with
 383 BWA-MEM. arXiv preprint arXiv:13033997. 2013.
- 384 18. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The
 385 Sequence Alignment/Map format and SAMtools. *Bioinformatics*. 2009;25
 386 16:2078-9. doi:10.1093/bioinformatics/btp352.
- 387 19. Visser I, Speekenbrink M: depmixS4: An R-package for hidden Markov
 388 models. *Journal of Statistical Software* 2010, 36(7):1-21.
- 389 20. Dobrynin P, Liu S, Tamazian G, Xiong Z, Yurchenko AA, Krasheninnikova K,
 390 Kliver S, Schmidt-Kuntzel A, Koepfli KP, Johnson W *et al*: Genomic legacy
 391 of the African cheetah, *Acinonyx jubatus*. *Genome Biol* 2015, 16:277.
- 392 21. Qiu Q, Zhang G, Ma T, Qian W, Wang J, Ye Z, et al. The yak genome and
 393 adaptation to life at high altitude. *Nat Genet*. 2012;44 8:946-9.
 394 doi:10.1038/ng.2343.
- 395 22. Wang K, Wang L, Lenstra JA, Jian J, Yang Y, Hu Q, et al. The genome
 396 sequence of the wisent (*Bison bonasus*). *Gigascience*. 2017;
 397 doi:10.1093/gigascience/gix016.
- 398 23. Bao W, Kojima KK and Kohany O. Repbase Update, a database of repetitive
 399 elements in eukaryotic genomes. *Mob DNA*. 2015;6:11. doi:10.1186/s13100-
 400 015-0041-9.
- 401 24. Tarailo-Graovac M and Chen N. Using RepeatMasker to identify repetitive
 402 elements in genomic sequences. *Curr Protoc Bioinformatics*. 2009;Chapter
 403 4:Unit 4 10. doi:10.1002/0471250953.bi0410s25.
- 404 25. Benson G. Tandem repeats finder: a program to analyze DNA sequences.
 405 *Nucleic Acids Res*. 1999;27 2:573-80.
- 406 26. Birney E, Clamp M and Durbin R. GeneWise and Genomewise. *Genome Res*.
 407 2004;14 5:988-95. doi:10.1101/gr.1865504.
- 408 27. Stanke M, Diekhans M, Baertsch R and Haussler D. Using native and
 409 syntenically mapped cDNA alignments to improve de novo gene finding.
 410 *Bioinformatics*. 2008;24 5:637-44. doi:10.1093/bioinformatics/btn013.
- 411 28. Blanco E, Parra G and Guigo R. Using geneid to identify genes. *Current*
 412 *protocols in bioinformatics*. 2007;Chapter 4:Unit 4.3.
 413 doi:10.1002/0471250953.bi0403s18.
- 414 29. Burge CB and Karlin S. Finding the genes in genomic DNA. *Curr Opin Struct*
 415 *Biol*. 1998;8 3:346-54.
- 416 30. Kielbasa SM, Wan R, Sato K, Horton P and Frith MC. Adaptive seeds tame
 417 genomic sequence comparison. *Genome Res*. 2011;21 3:487-93.
 418 doi:10.1101/gr.113985.110.
- 419 31. Yates A, Akanni W, Amode MR, Barrell D, Billis K, Carvalho-Silva D, et al.
 420 Ensembl 2016. *Nucleic Acids Research*. 2016;44 D1:D710-D6.

doi:10.1093/nar/gkv1157.

- 1 421
2 422 32. Dong Y, Xie M, Jiang Y, Xiao NQ, Du XY, Zhang WG, et al. Genomic data of
3 423 the domestic goat (*Capra hircus*). GigaScience Database
4 424 <http://dxdoiorg/105524/100082>. 2013.
5 425 33. Li L, Stoeckert CJ, Jr. and Roos DS. OrthoMCL: identification of ortholog
6 426 groups for eukaryotic genomes. *Genome Res.* 2003;13 9:2178-89.
7 427 doi:10.1101/gr.1224503.
8 428 34. De Bie T, Cristianini N, Demuth JP and Hahn MW. CAFE: a computational
9 429 tool for the study of gene family evolution. *Bioinformatics.* 2006;22 10:1269-
10 430 71. doi:10.1093/bioinformatics/btl097.
11 431 35. Loytynoja A and Goldman N. An algorithm for progressive multiple alignment
12 432 of sequences with insertions. *Proc Natl Acad Sci U S A.* 2005;102 30:10557-
13 433 62. doi:10.1073/pnas.0409137102.
14 434 36. Stamatakis A. RAxML version 8: a tool for phylogenetic analysis and post-
15 435 analysis of large phylogenies. *Bioinformatics.* 2014;30 9:1312-3.
16 436 doi:10.1093/bioinformatics/btu033.
17 437 37. Yang Z. PAML 4: phylogenetic analysis by maximum likelihood. *Mol Biol*
18 438 *Evol.* 2007;24 8:1586-91. doi:10.1093/molbev/msm088.
19 439 38. Benton MJ and Donoghue PC. Paleontological evidence to date the tree of life.
20 440 *Mol Biol Evol.* 2007;24 1:26-53. doi:10.1093/molbev/msl150.
21 441 39. Rambaut A and Drummond A. Tracer v1. 5 Available from [http://beast.bio.](http://beast.bio.ac.uk/Tracer)
22 442 [ac.uk/Tracer](http://beast.bio.ac.uk/Tracer). Accessed, 2013.
23 443 40. Li M, Tian S, Jin L, Zhou G, Li Y, Zhang Y, Wang T, Yeung CK, Chen L, Ma J
24 444 *et al*: Genomic analyses identify distinct patterns of selection in domesticated
25 445 pigs and Tibetan wild boars. *Nat Genet* 2013, 45(12):1431-1438.
26 446 41. Wang YX, Zheng YM: ROS-dependent signaling mechanisms for hypoxic
27 447 Ca(2+) responses in pulmonary artery myocytes. *Antioxid Redox Signal* 2010,
28 448 12(5):611-623.
29 449 42. Pijacka W, Moraes DJ, Ratcliffe LE, Nightingale AK, Hart EC, da Silva MP,
30 450 Machado BH, McBryde FD, Abdala AP, Ford AP: Purinergic receptors in the
31 451 carotid body as a new drug target for controlling hypertension. *Nature* 2016,
32 452 201:6.
33 453 43. Cai WK, Sakaguchi M, Kleinridders A, Gonzalez-Del Pino G, Dreyfuss JM,
34 454 O'Neill BT, Ramirez AK, Pan H, Winnay JN, Boucher J *et al*: Domain-
35 455 dependent effects of insulin and IGF-1 receptors on signalling and gene
36 456 expression. *Nat Commun* 2017, 8.
37 457 44. Rudovich N, Pivovarova O, Fisher E, Fischer-Rosinsky A, Spranger J, Mohlig
38 458 M, Schulze MB, Boeing H, Pfeiffer AF: Polymorphisms within insulin-
39 459 degrading enzyme (IDE) gene determine insulin metabolism and risk of type 2
40 460 diabetes. *J Mol Med (Berl)* 2009, 87(11):1145-1151.
41 461 45. Chung WK, Leibel RL: Molecular physiology of syndromic obesities in
42 462 humans. *Trends Endocrinol Metab* 2005, 16(6):267-272.
43 463 46. Charest-Marcotte A, Dufour CR, Wilson BJ, Tremblay AM, Eichner LJ, Arlow
44 464 DH, Mootha VK, Giguere V: The homeobox protein Prox1 is a negative

465 modulator of $ERR\alpha$ / $PGC-1\alpha$ bioenergetic functions. *Genes Dev* 2010,
466 24(6):537-542.

467 47. Li H, Durbin R: Inference of human population history from individual whole-
468 genome sequences. *Nature* 2011, 475(7357):493-496.

469 48. Wang Y, Song F, Zhu J, Zhang S, Yang Y, Chen T, et al. GSA: Genome
470 Sequence Archive. *Genom. Proteom. Bioinform.* 2017;15 1:14-8.
471 doi:10.1016/j.gpb.2017.01.001.

472 49. Members BIGDC. The BIG Data Center: from deposition to integration to
473 translation. *Nucleic Acids Res.* 2017;45 D1:D18-D24.
474 doi:10.1093/nar/gkw1060.

475 50. Yang Y, Wang Y, Zhao Y, Zhang X, Li R, Chen L et al. Supporting data for
476 "The genome of the Marco Polo Sheep (*Ovis ammon polii*)" *GigaScience*
477 Database. 2017. <http://dx.doi.org/10.5524/100366>

1 481 **Figure 1. Phylogenetic relationships and genomic comparisons between Marco**
2
3
4 482 **Polo Sheep and other mammals.** (a) Divergence time estimates for the nine mammals
5
6 483 generated using MCMCtree and the 4-fold degenerate sites. The red dots correspond to
7
8
9 484 calibration points and the divergence times. Divergence time estimates (Mya) are
10
11 485 indicated above the appropriate nodes; blue nodal bars indicate 95% confidence
12
13 486 intervals. Gene orthology was determined by comparing the genomes with the
14
15 487 OrthoMCL software. (b) A Venn diagram of the shared orthologues among Marco Polo
16
17 488 Sheep, sheep, goat, taurine cattle and horse. Each number represents a gene family
18
19
20 489 number. (c) The box plot shows the ratio of non-synonymous to synonymous mutations
21
22
23 490 (Ka/Ks) for Marco Polo Sheep, sheep, goat, taurine cattle, horse and human.
24
25
26
27

28 491

29
30 492
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

1 493 **Additional files**

2
3 494 **Figure S1.** 21-mer-based analysis carried out to estimate the size of the Marco Polo
4
5
6 495 Sheep genome.

7
8
9 496 **Figure S2.** GC content distribution for the genomes of Marco Polo Sheep, goat and
10
11 497 sheep.

12
13
14 498 **Figure S3.** FRCurve of five genome assemblies.

15
16
17 499 **Figure S4.** The distribution of observed heterozygosity stats within Marco Polo Sheep
18
19 500 genome.

20
21
22 501 **Figure S5.** Counts of InDels in coding regions, showing an enrichment of multiples
23
24
25 502 of three bases.

26
27
28 503 **Figure S6.** Comparison of gene structure characteristics with those of other
29
30 504 mammals.

31
32
33 505 **Figure S7.** Comparison of gene structure characteristics of the 1:1 orthologs in the
34
35
36 506 five mammals.

37
38
39 507 **Figure S8.** Comparison of the repeat content in the intron regions among Marco Polo
40
41
42 508 Sheep, Sheep (Oar_v3.1) and Goat (CHIR_1.0).

43
44
45 509 **Figure S9.** Summary of the number of chromosomes to which a given scaffold of the
46
47 510 Marco Polo Sheep genome could be aligned.

48
49
50 511 **Figure S10.** Divergence between Marco Polo Sheep, sheep and goat.

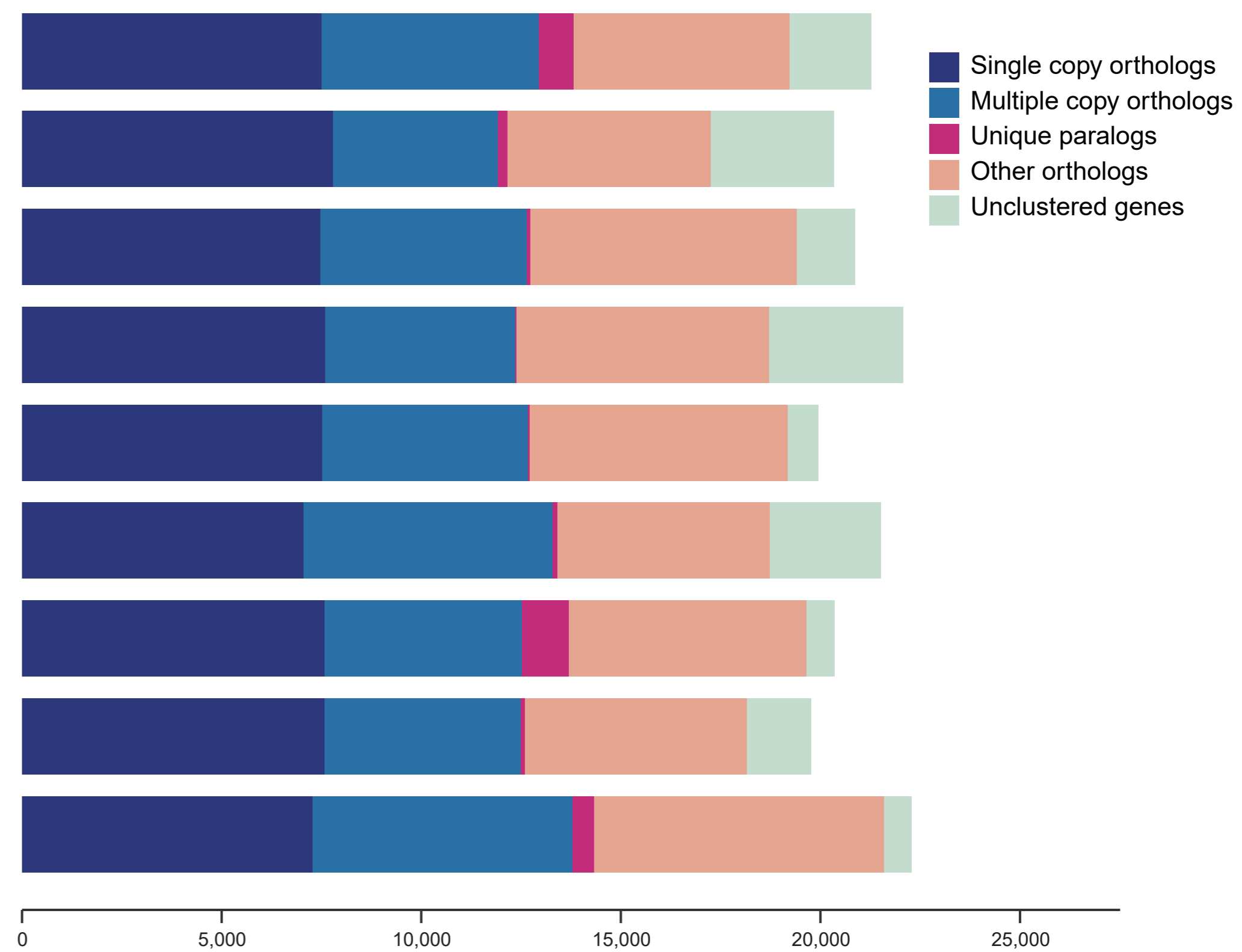
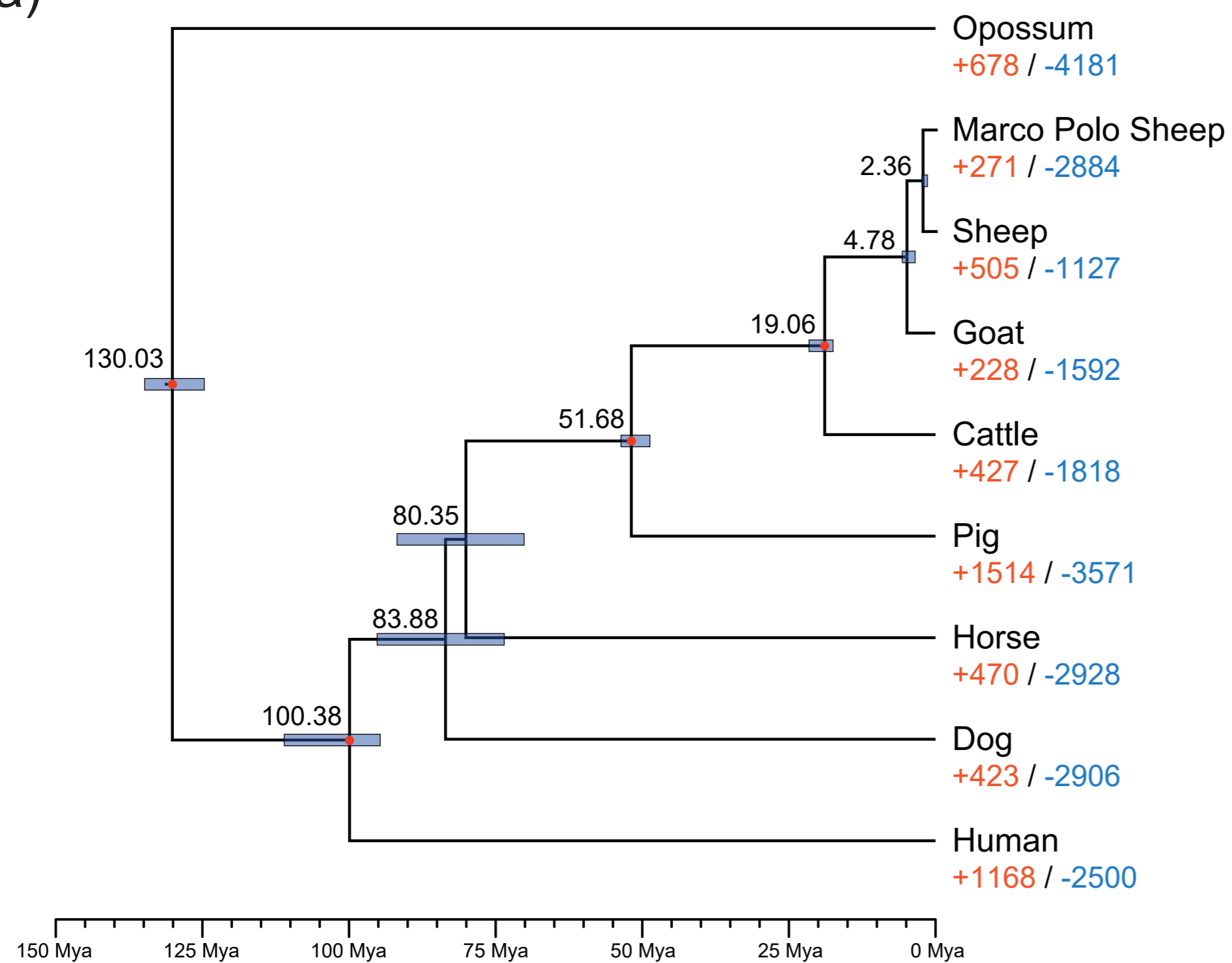
51
52
53 512 **Figure S11.** Synteny relationship between Marco Polo Sheep and sheep.

54
55
56 513 **Figure S12.** Synteny relationship between Marco Polo Sheep and goat.

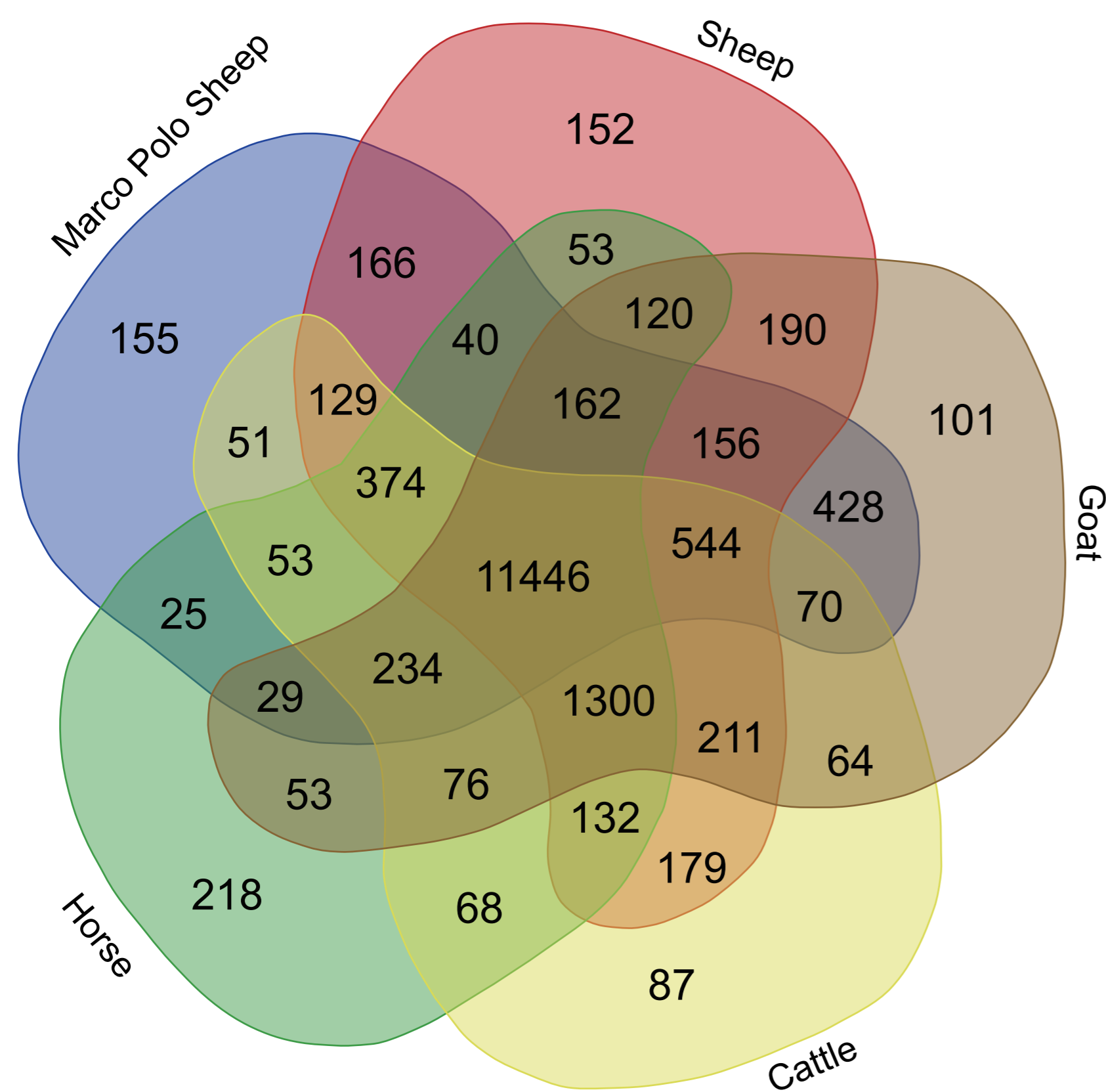
1 514 **Figure S13.** Density of breakpoints (number per million bases) in different regions of
2
3
4 515 the genome.
5
6 516 **Figure S14.** Phylogeny relationships between Marco Polo Sheep and other mammals.
7
8
9 517 **Figure S15.** Demographic history of Marco Polo Sheep.
10
11
12 518 **Table S1.** Summary of sequenced reads.
13
14
15 519 **Table S2.** Estimation of genome size based on 21-mer statistics.
16
17 520 **Table S3.** Statistics for the final assemblies of the Marco Polo Sheep genome.
18
19
20 521 **Table S4.** Numbers of reads mapped to the assembled Marco Polo Sheep genome.
21
22
23 522 **Table S5.** Summary of CEGMA analysis results.
24
25
26 523 **Table S6.** Summary of BUSCO analysis results obtained by counting matches to 4104
27
28 524 single-copy orthologs (mammalia_odb9).
29
30
31 525 **Table S7.** The distribution of SNVs in the Marco Polo Sheep genome.
32
33
34 526 **Table S8.** Genes located in the low heterozygosity regions.
35
36
37 527 **Table S9.** The distribution of InDels in the wisent genome.
38
39
40 528 **Table S10.** Prediction of repetitive elements in the assembled Marco Polo Sheep
41
42 529 genome.
43
44
45 530 **Table S11.** Classification of interspersed repeats in the assembled Marco Polo Sheep
46
47 531 genome.
48
49
50 532 **Table S12.** Data on all species used during the genome analysis.
51
52
53 533 **Table S13.** Prediction of protein-coding genes in the Marco Polo Sheep.
54
55
56 534 **Table S14.** Comparative gene statistics.
57
58
59 535 **Table S15.** Functional annotation of predicted genes in the Marco Polo Sheep.
60
61
62
63
64
65

1 536 **Table S16.** Summary statistics of non-coding RNAs in the Marco Polo Sheep.
2
3 537 **Table S17.** Summary of synteny alignments.
4
5
6 538 **Table S18.** Summary of breakpoints between Marco Polo Sheep, sheep and goat.
7
8
9 539 **Table S19.** Summary statistics of gene families in 9 species.
10
11 540 **Table S20.** GO enrichment analysis of the expanded gene families in the Marco Polo
12
13
14 541 Sheep lineage.
15
16
17 542 **Table S21.** Candidate positively selected genes (PSGs) in the Marco Polo Sheep
18
19
20 543 lineage.
21
22
23 544
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

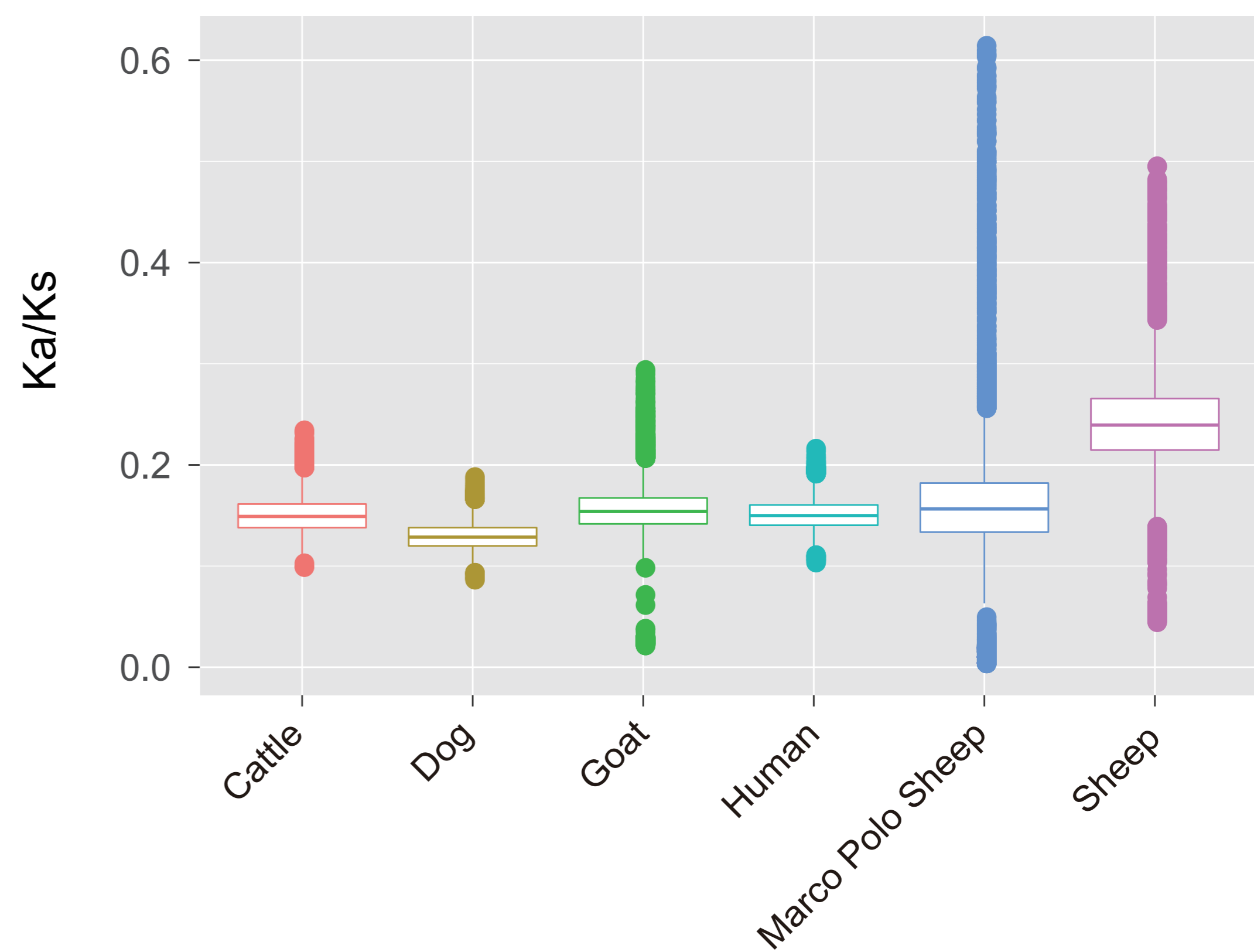
(a)



(b)



(c)

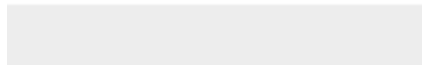




[Click here to access/download](#)

Supplementary Material

Macro polo sheep - Supplementary files.pdf



Dear Editor,

Thank you very much for returning our manuscript GIGA-D-17-00160 entitled “The genome of the Marco Polo Sheep (*Ovis ammon polii*)”.

We submit the revised manuscript here and hope the revised manuscript is more suitable for the publication in Giga Science. If you have any questions, please do not hesitate to contact the corresponding author at any time.

Thank you again for your time and efforts in handling our manuscript.

Best wishes,

Kun Wang

Center for Ecological and Environmental Sciences, Northwestern Polytechnical University, Xi’an 710072, China