

Manuscript Number:	GIGA-D-17-00185R1	
Full Title:	The draft genome sequence of forest musk deer (<i>Moschus berezovskii</i>)	
Article Type:	Data Note	
Funding Information:	National Key Program of Research and Development, Ministry of Science and Technology (2016YFC0503200)	Dr. Bisong Yue
	National Natural Science Foundation of China (31702032)	Dr. Wenhua Qi
Abstract:	<p>Background: The forest musk deer, <i>Moschus berezovskii</i>, is one of seven musk deer (<i>Moschus</i> spp.) and is distributed in Southwest China. Akin to other musk deer, the forest musk deer has been traditionally, and is currently, hunted for its musk (i.e. global perfume industry). Considerable hunting pressure and habitat loss has caused significant population declines and therefore the Chinese government commenced captive breeding programs for musk harvesting in the 1950s. However, the prevalence of fatal diseases is considerably restricting population increases. Disease severity and extent is exacerbated by inbreeding and genetic diversity declines in captive musk deer populations. It is essential for the physical and genetic health of captive and wild forest musk deer populations to improve knowledge of its immune system and genome. We have thus sequenced the whole genome of the forest musk deer, completed the genomic assembly and annotation, and performed preliminary bioinformatic analyses.</p> <p>Findings: A total of 407 Gb raw reads from whole-genome sequencing was generated by the Illumina HiSeq4000 platform. The final assembly genome is around 2.72 Gb, with a contig N50 length of 22.6 kb and a scaffold N50 length 2.85 Mb. We identified 24,352 genes, and found 42.05% of the genome is composed of repetitive elements. We also detected 1,236 olfactory receptor genes. The genome-wide phylogenetic tree indicated that the forest musk deer was within the order Artiodactyla, and it appeared as the sister clade of four members of family Bovidae. In total, 576 genes were under positive selection in the forest musk deer lineage.</p> <p>Conclusions: We provide the first genome sequence and gene annotation for the forest musk deer. The availability of these resources will be very useful for the conservation and captive breeding for this Endangered and economically important species, and for reconstructing the evolutionary history of the order Artiodactyla.</p>	
Corresponding Author:	Zhenxin Fan, Ph.D. Sichuan University Chengdu, Sichuan CHINA	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:	Sichuan University	
Corresponding Author's Secondary Institution:		
First Author:	Zhenxin Fan, Ph.D.	
First Author Secondary Information:		
Order of Authors:	Zhenxin Fan, Ph.D.	
	Bisong Yue	
	Wujiao Li	
	Chaochao Yan	
	Jing Li	

	Yongmei Shen
	Wenhua Qi
Order of Authors Secondary Information:	
Response to Reviewers:	<p>Response to the comments</p> <p>Reviewer #1</p> <p>Overall the genome reflects the industry standard for sequencing and assembly. The paper reads as a series of independent analyses that have not been integrated or justified. Some add virtually nothing to the paper, while others need to be flushed out a bit further (i.e. more detail and justification). If there is a word constraint I have highlighted two sections that could be dropped thereby providing room to expand on other sections. The english is generally good but a careful proof by a native speaker would be useful.</p> <p>1. Please provide the genome statistics after the short-read assembly, as this is useful information to the reader. Response: We added the information into the manuscript (Lines 125 - 132) and Table 2.</p> <p>2. Statements about length relative to other ungulates should be avoided, as this is dependent on the assembly quality and strategy. I do not find it biologically meaningful. Response: We deleted them.</p> <p>3. The TE analysis needs to be rethought; simply identifying repeats is not equivalent to identifying TEs. There is a severe lack of information provided here and it needs to be flushed out or dropped. How were the four types defined for example? Response: We explained within the manuscript (Lines 218 - 232).</p> <p>4. The Olfactory gene detection is not justified and it's not clear why this warrants it's own section. Response: We combined the olfactory receptor genes section into the "Gene families" section, and simplified it (Lines 259 - 263).</p> <p>5. The phylogenetic analysis does not add anything given there were no cervid genomes included; what was the point of this? Response: Because there are no Cervidae genomes, we could not fully address the phylogenetic position of forest musk deer at the complete genome sequences level in our manuscript. However, the phylogenetic tree still gave us the general picture of the phylogenetic relationship between forest musk deer and its related species, which their genomes are available now, and the results also indicated that we should add the genome sequences from family Cervidae in the further investigation. Therefore, we still keep the phylogenetic analysis in the manuscript.</p> <p>6. Again, no closely related species are included in the selection analysis. What was the point of this? Would you expect 500 genes to deviate from a closely related genome? Is this just chance? Response: We agreed with both the editor and reviewer, the selection analysis is not very useful here and also need more work for it. Therefore, we removed the selection analysis section, and more focused on the sequencing, assembly and annotation parts.</p>

Reviewer #2

General comments

In this article the authors present the first draft genome of the forest musk deer *Moschus berezovskii*. They provide a brief description of the sequencing, assembly and annotation process. This is a typical draft genome paper with little biological insight, but considering the status of the species and the little amount of available data on it, I believe that this contribution will be of great value to the community - if data and results are made available. I have two main concerns with this paper:

(1) I would need more details about the methods. It should be made possible to redo all the experiments and all the analyses that are mentioned in the manuscript. In particular, parameters and versions of each software tool have to be precised. Experimental protocols need more details too. One way to do this is to provide some more information in the main text and to complete with all the details in the Supplementary Methods. So far this document only contains the description of the phylogenetic analysis. The same should be done for the others.

Response: We added the parameters, software versions, and details of the protocols throughout the whole manuscript.

(2) The genomic sequence and its annotation should be made available. I could not find them. Sequencing reads have been deposited in the SRA but I would appreciate if the authors provide the results from the assembly (fasta format) and from the gene annotation (gff or gff3 format for instance). This is actually the main value of the study.

Response: We have uploaded the information to GigaScience's database, and we have communicated with the editor to make sure all the data is available.

More specific comments to the authors

1. Sampling/sequencing

L112-114: please provide more details about library construction (DNA extraction, protocol, kit, etc). I expect some to be PE and others mate-pairs. Could you precise? Also, I am not sure the read length is specified.

Response: We extracted genomic DNA from tissue samples using the Qiagen DNeasy Blood and Tissue Kit (Qiagen, Valencia, USA) following the manufacturer's protocol. We added the information into the "Sample information and sequencing" section (Lines 112 - 119). We also added the read length information into the Table 1. The details were given in the new supplementary note.

L115-116: what kind of filtering/cleaning has been made and how? Please provide details about quality and adapter trimming, including the name of the software.

Response: We provided the details within the supplementary note (sequencing and filtering raw data section).

2. Assembly

L120-121: how did you estimate the genome size? Method, software, version? Why 17-mers?

Response: We used GCE (v1.0) to perform kmer analysis, and the 17-mer was suggested by the GCE results. The details could be found within the supplementary note.

L122: "the assembly was first analyzed by SOAPdenovo2" => don't you mean "generated"? The assembly needs to be produced before being analyzed.

Response: Yes, it should be "generated" and we corrected.

L125: what was the proportion of gaps before and after gapcloser?

Response: We explained it within the supplementary note.

L124: how was SSPACE used? How many scaffolds before and after? Also, please precise the version.

Response: We explained it within the supplementary note.

L129: Table 2 is way too short -only three numbers!- to give a decent description of the

assembly. Please give the number of contigs, of scaffolds, the size of the longest ones, the GC%, etc...

See the same tables in similar publications from the same journal, for instance:

<https://academic.oup.com/gigascience/article/4077042/The-draft-genome-sequence-of-a-desert-tree-Populus>

<https://academic.oup.com/gigascience/article/6/8/1/4004833/Draft-genome-of-the-Antarctic-dragonfish>

<https://academic.oup.com/gigascience/article/6/6/1/3748232/Draft-genome-of-the-lined-seahorse-Hippocampus>

More generally, please also consider these previous publications to get an idea about the amount of details that are expected from this kind of report.

Response: We added more information within Table 2.

L132: CEGMA + BUSCO: Cegma is no longer maintained and should not be used anymore.

Response: Yes, we deleted the Cegma result and former Table S1 (the results for Cegma).

L136: please cite the study that generated the RNA-seq data you used.

Response: We added the SRA ID for the RNA-seq data (SRR2098995 and SRR2098996; Line 138) since there was no publication.

L138: the proportion of mapped reads is high and this is a good point, but it would be more informative to also show the proportion of concordant pairs, assuming that the RNA-seq data is PE. This assembly section is rather descriptive and technical but it is difficult to estimate up to which point all these steps were useful. A nice way to show the value of this work could be to compare the number of mapped reads and concordant read pairs from the RNA-seq and DNA-seq libraries (of the different sizes) before and after the gap filling and scaffolding part.

Response: In total of 92.73% PE reads were concordantly aligned to forest musk deer genome. We added the information into the manuscript (Lines 137 - 138).

3. Annotation

L147: how was Augustus trained? Was a training set of known genes provided to estimate the parameters?

Response: Yes, we used some protein sequences of *Bos taurus* to train the model for Augustus (version 3.0.3). We provided the information within supplementary note.

L148: "analyzed" => aligned, I guess

Response: Yes, we corrected the word.

L150-151: what is "software solar"? Please explain how GeneWise was used and provide details about potential filtering and other steps after the blast.

Response: We explained within the main text (Lines 178 – 180)

L153: Trinity has been used in both genome guided and de novo mode: why? What were the differences? Please provide the parameters. How did you merge the results?

Response: This method was followed by Trinity manual in the section of ' Build a Comprehensive Transcriptome Database Using Genome-guided and De novo RNA-Seq Assembly ', website is <https://pasapipeline.github.io/>.

L154: Please cite EVM properly. How did you use it? What did you choose for the confidence weights? Could you provide the input files from the distinct sources before merging?

Response: We cited the reference "Automated eukaryotic gene structure annotation using EVIDENCEModeler and the Program to Assemble Spliced Alignments, Genome Biology 2008". Other details could be found in the supplementary note.

L156: manual annotation can be a huge amount of work. Were there many modifications made? If so, it would be interesting and probably impressive to illustrate the contribution of this work by comparing annotation stats (see below) or other metrics before and after this polishing step, and/or to describe the most common corrections (gene splitting/merging? Splice site fixing? etc). That is only a suggestion.

	<p>Response: Yes, the manual annotation will require lots of work. Therefore, we only manually checked the scaffolds that size longer than 1Mb. In addition, we focused on gene splitting or merging.</p> <p>Before the GO functional annotation, something that is missing is the description of the annotation with more statistics than just the number of genes, especially given the draft status of the genome assembly. In particular a simple table could present some stats about the gene length distribution (min, max, median, average), distribution of predicted ORFs/CDS (idem), number of exons per gene (idem: min, max, mean and median).</p> <p>Response: We added a new supplementary table (current Table S2) to provide the length of the gene, CDS and exon.</p> <p>Also, it would be useful to illustrate the quality of the provided annotation by comparing it with other available datasets. For instance, what is the percentage of RNA-seq mapped reads that fall within annotated exons? That are consistent with the predicted gene models? How do the annotated transcripts compare with those from the already published transcriptome study (ies)?</p> <p>Response: Furthermore, we downloaded musk gland RNA-seq data (SRA accession: SRR2098995 and SRR2098996) of forest musk deer from NCBI to evaluate the assembly. We found that 99.3% of the total PE reads could be aligned (92.73% aligned concordantly) to the assembled forest musk deer genome by Bowtie2 (version 2.2.5). We added this within the main text (Lines 133 - 167).</p> <p>4. Repeats L190: "We also analyzed the degree of divergence for each type of TE" => How? Again: method, software, version, and parameters. Same for MSDB. Response: We updated the method, software, version and parameters (Lines 219 – 239).</p> <p>The number of TEs is compared across species: could the authors make sure that the same method was used to detect them in each species? Otherwise it could just be due to the method. Please keep in mind that all these annotated "genes", "TEs" etc, are only computational predictions. Response: Yes, they were detected with different methods although the methods were very similar. We still should be careful when we compared the numbers.</p> <p>5. Gene families See general concern (1). Please also indicate the version of the ENSEMBL and NCBI annotation. Response: Yes, we added the information into the table (current Table S5).</p> <p>6. Olfactory Receptor genes L218: typos in "pseudogenes" and "truncated". Response: Thanks, we corrected the words.</p> <p>L219: The number of OR genes is compared across species and a "degeneration of OR genes in primates" is mentioned. Couldn't the difference be due to the fact that these OR genes were not annotated the same way between species (Sup Tables)? Response: Yes, therefore, we removed the comparison, and reviewer 1 also asked us to simplify this part (Lines 260 - 264).</p> <p>Table S2: $19 / 303 = 20.51\%$? ("missing busco" part) Response: It was a mistake, we corrected the numbers and now, it is now within Table S1.</p> <p>Table S5: It is ensembl, not ensemble. Please precise the version. Response: We corrected the word, and also provided the version and other related information.</p>
Additional Information:	
Question	Response
Are you submitting this manuscript to a	No

special series or article collection?	
<p>Experimental design and statistics</p> <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>	Yes
<p>Resources</p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite 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>	Yes
<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>	Yes

1 **The draft genome sequence of forest musk deer (*Moschus berezovskii*)**

2
3
4
5
6
7
8

9 Zhenxin Fan^{1,†}, Wujiao Li^{1,†}, Jiazheng Jin¹, Kai Cui¹, Chaochao Yan¹, Changjun

10 Peng¹, Zuoyi Jian¹, Megan Price¹, Xiuyue Zhang¹, Yongmei Shen³, Jing Li¹, Wenhua

11 Qi^{2,*}, Bisong Yue^{1,*}

12
13
14
15
16

17 ¹Key Laboratory of Bioresources and Ecoenvironment (Ministry of Education),

18 College of Life Sciences, Sichuan University, Chengdu 610064, People's Republic

19 of China

20 ²College of Life Science and Engineering, Chongqing Three Gorges University,

21 Chongqing 404100, People's Republic of China

22 ³Sichuan Engineering Research Center for Medicinal Animals, Xichang 615000,

23 People's Republic of China

24 † Contributed equally to this work

25 * **Corresponding author:** Bisong Yue (bsyue@scu.edu.cn), and Wenhua Qi

26 (wenhuaqi357@163.com)

27 **Emails:** zxfan@scu.edu.cn (Zhenxin Fan); hnnd059@gmail.com (Wujiao Li);

28 wenhuaqi357@163.com (Wenhua Qi); jinjiazhengxiao@163.com (Jiazheng Jin);

29 2015222040118@stu.scu.edu.cn (Kai Cui); yccvican@gmail.com (Chaochao Yan);

30 jj-5380682@163.com (Changjun Peng); jzuoyi@126.com (Zuoyi Jian);

31 meganprice@scu.edu.cn (Megan Price); zhangxy317@126.com (Xiuyue Zhang);

32 Yongmei Shen (810316122@qq.com); Jing Li (ljtjf@126.com); bsyue@scu.edu.cn

33 (Bisong Yue)

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

24 **Abstract**

25 **Background:** The forest musk deer, *Moschus berezovskii*, is one of seven musk
26 deer (*Moschus* spp.) and is distributed in Southwest China. Akin to other musk
27 deer, the forest musk deer has been traditionally, and is currently, hunted for its
28 musk (i.e. global perfume industry). Considerable hunting pressure and habitat
29 loss has caused significant population declines and therefore the Chinese
30 government commenced captive breeding programs for musk harvesting in the
31 1950s. However, the prevalence of fatal diseases is considerably restricting
32 population increases. Disease severity and extent is exacerbated by inbreeding
33 and genetic diversity declines in captive musk deer populations. It is essential for
34 the physical and genetic health of captive and wild forest musk deer populations
35 to improve knowledge of its immune system and genome. We have thus
36 sequenced the whole genome of the forest musk deer, completed the genomic
37 assembly and annotation, and performed preliminary bioinformatic analyses.

38 **Findings:** A total of 407 Gb raw reads from whole-genome sequencing was
39 generated by the Illumina HiSeq4000 platform. The final assembly genome is
40 around 2.72 Gb, with a contig N50 length of 22.6 kb and a scaffold N50 length
41 2.85 Mb. We identified 24,352 genes, and found 42.05% of the genome is
42 composed of repetitive elements. We also detected 1,236 olfactory receptor
43 genes. The genome-wide phylogenetic tree indicated that the forest musk deer
44 was within the order Artiodactyla, and it appeared as the sister clade of four
45 members of family Bovidae. In total, 576 genes were under positive selection in
46 the forest musk deer lineage.

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

47 **Conclusions:** We provide the first genome sequence and gene annotation for the
48 forest musk deer. The availability of these resources will be very useful for the
49 conservation and captive breeding for this Endangered and economically
50 important species, and for reconstructing the evolutionary history of the order
51 Artiodactyla.

52

53 **Keywords:** Forest musk deer; whole genome sequencing; genome assembly;
54 annotation; phylogeny

55

56

57

58

59

60

61

62

63

64

65

66

67

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

68 Data Description

69 1) Background

70 The seven musk deer species of genus *Moschus* are endemic to Asia, are
71 currently listed under Appendix II in CITES and are listed under Category I of the
72 State Key Protected Wildlife List of China [1-3]. All musk deer species are
73 considered as globally threatened, with six being listed as Endangered and one
74 as Vulnerable by the IUCN [4]. *Moschus* is the only extant genus of Moschidae and
75 musk deer are considered as primitive deer. The genus of musk deer is
76 characterized by the musk secreted by the scent glands of adult males [5]. The
77 forest musk deer (*Moschus berezovskii*) is one of the five recognized musk deer
78 species of China and have historically been distributed in Southwest China [6,7].
79 The forest musk deer has been listed as globally endangered, as Critically
80 Endangered on the 2015 China Red List, and is also on the State Key Protected
81 Wildlife List of China [4].

82 Musk deer have been hunted for thousands of years, as the musk has been
83 widely used in traditional Chinese medicines. In the last two centuries, hunting
84 of all musk deer species significantly increased for the global trade of the
85 commercially valuable musk secretion as an essential basis for perfume
86 manufacture [5]. Since the 1950s, populations of forest musk deer have declined
87 dramatically from poaching of deer for the musk pods (i.e. entire gland) and
88 significant habitat destruction [3,6,8]. As a consequence, the Chinese
89 government has encouraged musk-using enterprises to participate in artificial
90 breeding programs since the early 1950s [9]. The musk can be collected from

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

91 male musk deer in these captive populations without harvesting individuals,
92 further enhancing the commercial and conservation value of captive populations.

93 The captive population of the forest musk deer is the largest among all
94 the musk deer species [2,10]. The Miyaluo farming population in Sichuan
95 Province (China) was one of the earliest established captive breeding
96 populations. This population had grown rapidly to approximately 400 in 2010
97 [10]. However, the prevalence of fatal diseases is considerably restricting
98 population increases [11]. Common diseases of forest musk deer in the Miyaluo
99 population are dyspepsia, pneumonia, metritis, urinary stones and abscesses,
100 with abscesses being one of the most prevalent causes of death [7]. Disease
101 severity and extent is exacerbated by inbreeding and genetic diversity declines
102 in this and other captive musk deer populations.

103 It is essential for the physical and genetic health of captive and wild forest
104 musk deer populations to improve knowledge of its immune system and genome.
105 We have thus sequenced the whole genome of the forest musk deer,
106 subsequently completed the genomic assembly and annotation, and performed
107 preliminary bioinformatic analyses, such as phylogenetic tree, selection and gene
108 enrichments.

109

110 2) Sample information and sequencing

111 The thigh muscle sample was collected from a Miyaluo male forest musk deer
112 that naturally died (Sichuan Province, China) in 2015. We extracted genomic
113 DNA from the muscle sample using the Qiagen DNeasy Blood and Tissue Kit
114 (Qiagen, Valencia, USA) following the manufacturer's protocol. We constructed

115 six different insert size libraries: 230bp, 500bp, 2kb, 5kb, 10kb, and 15kb. These
116 libraries were sequenced by Illumina Hiseq 4000 platform at Novogene (Beijing,
117 China). A total of 407Gb of raw data were generated, after filtering out low
118 quality, duplicate and adaptor polluted reads. Approximately 360Gb of high-
119 quality reads were retained for genome assembly (Table 1).

120

121 3) Genome assembly and evaluation

122 We use GCE (version 1.0) to performed k-mer (17-mer) analysis by short insert
123 size library reads before assembly, and the forest musk deer genome size was
124 estimated to be 2.95Gb (Figure S1). The assembly was first generated by
125 SOAPdenovo2 [12] with the parameters set as “all -d 2 -M 2 -k 35”. Intra-
126 scaffold gaps were filled using Gapcloser (version 1.12) with reads from 230bp
127 and 500bp libraries, and then SSPACE (version 3.0) [13] was used to build super-
128 scaffolds. After scaffolding by SSPACE, we used Gapcloser again to fill gaps.
129 Finally we obtained the forest musk deer genome with a size of 2.72Gb (all the
130 sequences with length shorter than 300bp were removed) with 125.7Mb gap
131 sequences unsolved. The N50s of contigs and scaffolds of forest musk deer
132 genome were 22.6kb and 2.85Mb, respectively (Table 2).

133 We used BUSCO (version 3.0) to evaluate the genome complement.
134 BUSCO results showed that 84.5% of the eukaryotic single-copy genes were
135 captured (Table S1). Furthermore, we downloaded musk gland RNA-seq data
136 (SRA accession: SRR2098995 and SRR2098996) of forest musk deer from NCBI
137 to evaluate the assembly. We found that 99.3% of the total PE reads could be
138 aligned (92.73% aligned concordantly) to the assembled forest musk deer

139 genome by Bowtie2 (version 2.2.5) [14]. These results showed our forest musk
140 deer genome was of high quality, and was suitable as a reference genome for the
141 family Moschidae.

142

143 4) Annotation

144 We combined the *de novo*, homology-based and transcriptome-based prediction
145 to identify protein-coding genes in the forest musk deer genome. The software
146 Augustus (version 3.2.1) [15] was used for *de novo* prediction based on the
147 parameter trained for forest musk deer. For homology prediction, protein
148 sequences from four mammals (human, pig, sheep and cattle) were analyzed
149 with TBLASTN (BLAST version 2.2.26) against forest musk deer genome.

150 Potential gene regions were joined by SOLAR (version 0.9.6) [16], and the coding
151 sequence with 500bp flanking sequence were cut down and re-aligned by
152 GeneWise (version 2.4.1 with parameters “- sum - genesf -gff”) [17]. For
153 transcriptome-based prediction, musk gland RNA-seq data were assembled by
154 Trinity with genome guide and *de novo* mode, respectively. The gene structures
155 were obtained by PASA pipeline (version 2.0.2) [18]. We used EVM (version 1.1.1)
156 to integrate the above evidence and obtained a consensus gene set [19]. Apollo
157 (version 1.11.6) was performed to manually inspect gene structure in scaffolds
158 of sizes above 1Mb to gain a more accurate gene structure. We consequently
159 found a total of 24,352 genes predicted to be present in the forest musk deer
160 genome. We also provided the length of genes in Table S2.

161 Functional annotation of forest musk deer genes was undertaken based
162 on the best match derived from the alignments to proteins annotated in Swiss-

1 163 Prot and TrEMBL databases [20]. Functional annotation used BlastP tools with
2 164 the same E-value cut-off of 1E-5. We also annotated proteins against the NCBI
3
4 165 non-redundant (nr) protein database. The outputs of blast searching against the
5
6 166 NCBI nr protein database were imported into BLAST2GO (B2G4PIPE v2.5) for
7
8 167 Gene Ontology (GO) [21] term mapping. Term mapping used annotated motifs
9
10 168 and domains using InterProScan (interproscan-5.18-57.0) [22] by searching
11
12 169 against publicly available databases. To find the best match for each gene, KEGG
13
14 170 pathway maps were used by searching KEGG databases [23] through the KEGG
15
16 171 Automatic Annotation Server (KAAS) using the bi-directional best hit (BBH)
17
18 172 method. In total, 23,023 out of 24,352 (94.5%) protein-coding genes were
19
20 173 searched within the publicly available functional databases of TrEMBL, Swiss-
21
22 174 Prot, Interpro, GO and KEGG. Of which, 22,696 (93.20% TrEMBL), 18,771 (77.08%
23
24 175 Swiss-Prot), 22,221 (91.12% Interpro), 15,736 (64.62% GO) and 10,846 (44.54%
25
26 176 KEGG) genes showed significant similarity matches (Figure 1; Table 3). The
27
28 177 functional comparisons with two closely related species (cattle and sheep) for
29
30 178 GO classification were submitted to the WEGO [24] (Figure S2).
31
32
33
34
35
36
37
38
39
40
41
42

43 180 5) Repetitive sequences and transposable elements

44
45 181 Transposable elements (TEs) and other repeats make up a substantial fraction of
46
47 182 mammalian genomes and contribute to gene and/or genome evolution [25]. The
48
49 183 TE content, type, copy number, subfamily, and divergence rate were investigated
50
51 184 in the forest musk deer genome based on two strategies: the library based
52
53 185 strategy of RepeatMasker [26] and the *de novo* based strategy of RepeatScout
54
55 186 [27]. The forest musk deer genome has large numbers of TEs, comprising 42.05%
56
57
58
59
60
61
62
63
64
65

187 of the genome (Table S3), which is similar to those of cattle (46.5%) [25] and
188 goats (42.2%) [28]. The 23 different types of TEs have been grouped for the four
189 different types of TEs, such as DNA transposons, LTR, LINE, and SINE
190 retrotransposons (Figure S3). The LINEs were the most common repeats in forest
191 musk deer genome; followed by SINEs > LTR > DNA. We also analyzed the
192 degree of divergence for each type of TE in the forest musk deer genome. We
193 found there was a recent burst activity involving LINE transposons and a second,
194 older burst activity of LTR and DNA transposons (Figure S3).

195 A total of 542,135 microsatellites (simple sequence repeats, SSRs) were
196 identified by software MSDB [29] in the forest musk deer genome assembly
197 (Table S4), which accounted for 0.45% of its whole genome length.

198 Mononucleotide SSRs were the most abundant category, accounting for 41.75%
199 of all of the SSRs; followed by di- > tri- > tetra- > penta- > hexa
200 nucleotide SSRs (Table S4).

201

202 6) Gene families

203 To estimate species-specific and shared genes in the forest musk deer in
204 comparison to ten mammal species, we used orthoMCL [30] to define the
205 orthologous genes. We downloaded the genomes and gene annotations of the ten
206 additional species (human, horse, dog, cattle, mouse, yak, sheep, Tibetan
207 antelope, alpaca, and pig) from Ensembl [31] or NCBI (Table S5). In total, we
208 identified 18,855 homologous gene families shared by forest musk deer and the
209 ten additional species, 221 gene families that were specific to forest musk deer,
210 and 2,003 gene families found in the ten additional species but not in the forest

1 211 musk deer (Figure S4). In addition, we found 5,372 one-to-one orthologous
2 212 genes within forest musk deer and the ten species, which was used in
3
4 213 phylogenetic analyses. In addition, we detected olfactory receptor (OR) genes in
5
6 214 the forest musk deer genome by orfam (<https://github.com/jianzuoyi/orfam>)
7
8 215 since they formed the largest gene family in mammalian genomes [32]. In total,
9
10 216 we identified 1,236 OR genes, which included 866 intact, 266 pseduogenes, and
11
12 217 104 truncated genes.
13
14
15
16
17
18 218

219 7) Phylogenetic analysis

220 We constructed the phylogenetic trees based on Bayesian inference (BI) [33] and
221 maximum likelihood (ML) [34,35] analyses with the discovered 5,372 one-to-
222 one orthologous genes (Supplementary notes). All the different methods
223 generated the same topology and obtained the well-supported phylogenetic tree
224 (Figure 2). The forest musk deer was within the suborder Ruminantia, order
225 Artiodactyla, and it appeared as the sister clade of four members of family
226 Bovidae (sheep, yak, cattle, and Tibetan antelope). Since we do not have high
227 quality genome sequences for species within family Cervidae, the relationship
228 between Moschidae, Cervidae, and Bovidae at the genomic level is tentative and
229 needs further investigation.

230

231 **Conclusions**

232 Here, we report the first draft genome assembly of the forest musk deer genome,
233 a species that is of particular importance to China's ecology, biodiversity
234 conservation, economy, and medicine. The availability of the genome and these

1 235 results will be very useful for the conservation and captive breeding of this
2 236 Endangered and economically important species, and for reconstructing the
3
4 237 evolutionary history of the order Artiodactyla.
5
6

7
8 238

9
10
11 239 **Funding**

12
13
14 240 This work was supported by National Key Program of Research and
15
16 241 Development, Ministry of Science and Technology (2016YFC0503200), and
17
18 242 National Natural Science Foundation of China (31702032).
19
20

21
22 243

23
24
25 244 **Availability of supporting data**

26
27
28 245 The DNA sequencing data have been deposited into the NCBI Sequence Read
29
30 246 Archive (SRA) under the ID PRJNA317652.
31
32

33
34 247

35
36 248 **Conflicts of interest**

37
38
39 249 The authors declare that they have no competing interests.
40
41

42
43 250

44
45 251 **Author's contributions**

46
47
48 252 Z.F., X.Z., J.L., and B.Y. designed and supervised the project. Z.F., W.L., C.Y., J.J., C.P.,
49
50 253 J.Y., Y.S., and K.C. performed the bioinformatics analyses. M.P. revised the
51
52 254 manuscript. Z.F. and B.Y. wrote the manuscript.
53
54

55
56 255

57
58
59 256
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

257 **Figure Legend**

258 **Figure 1 Functional annotation statistics.** Venn diagram illustrating
259 distribution of high-score matches of the functional annotation in forest musk
260 deer genome from five public databases.

261 **Figure 2 Genome wide phylogenetic trees.** We constructed the phylogenetic
262 trees based on Bayesian inference and maximum likelihood analyses with 5,372
263 one-to-one orthologous genes between the forest musk deer and ten other
264 species.

265 **Figure S1 K-mer (k=17) distributions in forest musk deer genome.**

266 **Figure S2 GO comparative analysis and functional classification between**
267 **forest musk deer, sheep and cattle.**

268 **Figure S3 Distribution of divergence of each type of TEs in the forest musk**
269 **deer genome.** The divergence rate was calculated between the identified TE
270 elements in the genome and the consensus sequence in the TE library used.
271 SINEs: Short interspersed elements, LINEs: Long interspersed elements, LTR:
272 Long terminal repeat retrotransposon.

273 **Figure S4 Protein orthology comparison between different genomes.** There
274 were forest musk deer (*Moschus bweezovskii*), cattle (*Bos taurus*), yak (*Bos*
275 *grunniens*), sheep (*Ovis aries*), Tibetan antelope (*Pantholops hodgsonii*), alpaca
276 (*Vicugna pacos*), and pig (*Sus scrofa*), which representing Artiodactyla; human
277 (*Homo sapiens*, Primates), horse (*Equus caballus*, Perissodactyla), and dog (*Canis*
278 *lupus familiaris*, Carnivora), mouse (*Mus musculus*, Rodentia). For each animal,
279 proteins are represented by bars and classified according to orthoMCL analysis:
280 Single_copy (Oliver) include the common orthologs with the same number of

1 281 copies in different species; Multi_copy (Red) include the common orthologs with
2 282 different copy numbers in the different species; Unique (Magenta) include the
3
4 283 orthologs just in one species; Unclustered gene (Yellow) include the genes that
5
6 284 cannot be clustered into known gene families; Other (Blue) include the genes
7
8 285 that can be clustered into known gene families, but it not belongs to Single, Multi
9
10
11
12 286 or Unique.
13

14
15 287
16

17
18 288
19

20
21 289
22

23
24 290
25

26
27 291
28

29
30 292
31

32
33 293
34

35
36 294
37

38
39 295
40

41
42 296
43

44
45 297
46

47
48 298
49

50
51 299
52

53
54 300
55

56
57 301
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

302 **References**

- 303 1. Cobert GB, Hill JE. The mammals of the Indomalaysia region: a systematic
304 review. London (UK): Natural History Museum Publications Oxford
305 University Press. 1992.
- 306 2. Wu J, Wang W. The Musk Deer of China. The China Forestry Publishing House,
307 Beijing. 2006.
- 308 3. Sheng H, Liu Z. The Musk Deer in China. The Shanghai Scientific & Technical
309 Publishers, Shanghai. 2007.
- 310 4. IUCN. The IUCN Red List of Threatened Species. 2017
311 (<http://www.iucnredlist.org/>)
- 312 5. Green MJB. The distribution, status and conservation of the Himalayan musk
313 deer (*Moschus chrysogaster*). *Biol. Conserv.* 1986;35, 347–75.
- 314 6. Sheng H. Genus *Moschus* in China. In: Wang, S. (Ed.), *China Red Data Book of*
315 *Endangered Animals*. Science Press, Beijing. 1998.
- 316 7. Zhao K, Liu Y, Zhang X, et al. Detection and characterization of antibiotic-
317 resistance genes in *Arcanobacterium pyogenes* strains from abscesses of
318 forest musk deer. *J Med Microbiol.* 2011;60:1820-6.
- 319 8. Yang Q, Meng X, Xia L, Feng Z. Conservation status and causes of decline of
320 musk deer (*Moschus spp.*) in China. *Biol. Conserv.* 2003;109, 333–342.
- 321 9. Peng H, Liu S, Zou F, Zeng B, Yue B. Genetic diversity of captive forest musk
322 deer (*Moschus berezovskii*) inferred from the mitochondrial DNA control
323 region. *Anim Genet.* 2009;40(1):65-72.

- 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
- 324 10. Huang J, Li Y, Li P, et al. Genetic quality of the Miyaluo captive forest musk
325 deer (*Moschus berezovskii*) population as assessed by microsatellite loci.
326 *Biochemical Systematics & Ecology*, 2013;47(8):25-30.
 - 327 11. Lu X, Qiao J, Wu X, Su L. A review of mainly affected on musk-deer diseases:
328 purulent, respiratory system and parasitic diseases. *J Economic Anim*
329 2009, 13, 104–107.
 - 330 12. Luo R, Liu B, Xie Y, et al. SOAPdenovo2: an empirically improved memory-
331 efficient short-read de novo assembler. *Gigascience* 2012;1(1):18.
 - 332 13. Boetzer M, Henkel CV, Jansen HJ, Butler D, Pirovano W. Scaffolding pre-
333 assembled contigs using SSPACE. *Bioinformatics*. 2011;27(4):578-9.
 - 334 14. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nat*
335 *Methods*. 2012;9(4):357-9.
 - 336 15. Stanke M, Diekhans M, Baertsch R, Haussler D. Using native and syntenically
337 mapped cDNA alignments to improve de novo gene finding.
338 *Bioinformatics*. 2008;24(5):637-44.
 - 339 16. Yu XJ, Zheng HK, Wang J, Wang W, Su B. Detecting lineage-specific adaptive
340 evolution of brain-expressed genes in human using rhesus macaque as
341 outgroup. *Genomics*. 2006;88(6):745-751.
 - 342 17. Birney E, Clamp M, Durbin R. GeneWise and Genomewise. *Genome Res*.
343 2004;14(5):988-95.
 - 344 18. Haas BJ, Delcher AL, Mount SM, et al. Improving the Arabidopsis genome
345 annotation using maximal transcript alignment assemblies. *Nucleic Acids*
346 *Res*. 2003;31(19):5654-66.

- 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
- 347 19. Haas BJ, Salzberg SL, Zhu W, Pertea M, Allen JE, Orvis J, White O, Buell CR,
348 Wortman JR. Automated eukaryotic gene structure annotation using
349 EVIDENCEModeler and the Program to Assemble Spliced Alignments.
350 Genome Biol. 2008;9(1):R7.
- 351 20. Boeckmann B, Bairoch A, Apweiler R, et al. The SWISS-PROT protein
352 knowledgebase and its supplement TrEMBL in 2003. Nucleic Acids Res.
353 2003;31(1):365-70.
- 354 21. Gene Ontology Consortium. Gene Ontology annotations and resources.
355 Nucleic Acids Res. 2013;41(Database issue):D530-5.
- 356 22. Hunter S, Apweiler R, Attwood TK, et al. InterPro: the integrative protein
357 signature database. Nucleic Acids Res. 2009;37(Database issue):D211-5.
- 358 23. Kanehisa M, Goto S. KEGG: kyoto encyclopedia of genes and genomes. Nucleic
359 Acids Res. 2000;28(1):27-30.
- 360 24. Ye J, Fang L, Zheng H, et al. WEGO: a web tool for plotting GO annotations.
361 Nucleic Acids Res. 2006;34(Web Server issue):W293-7.
- 362 25. Adelson DL, Raison JM, Edgar RC. Characterization and distribution of
363 retrotransposons and simple sequence repeats in the bovine genome.
364 Proc Natl Acad Sci U S A. 2009;106(31):12855-60.
- 365 26. Smit AFA, Hubley R, Green P. 2016. RepeatMasker website and
366 server[CP/OL]. [2016-9-12][2016-10-15].
367 <http://www.repeatmasker.org/>.
- 368 27. Price AL, Jones NC, Pevzner PA. De novo identification of repeat families in
369 large genomes. Bioinformatics. 2005;21 Suppl 1:i351-8.

- 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
- 370 28. Dong Y, Xie M, Jiang Y, et al. Sequencing and automated whole-genome
371 optical mapping of the genome of a domestic goat (*Capra hircus*). *Nat*
372 *Biotechnol.* 2013;31(2):135-41.
- 373 29. Du L, Li Y, Zhang X, Yue B. MSDB: a user-friendly program for reporting
374 distribution and building databases of microsatellites from genome
375 sequences. *J Hered.* 2013;104(1):154-7.
- 376 30. Li L, Stoeckert CJ Jr, Roos DS. OrthoMCL: identification of ortholog groups for
377 eukaryotic genomes. *Genome Res.* 2003;13(9):2178-89.
- 378 31. Yates A, Akanni W, Amode MR et al. Ensembl 2016. *Nucleic Acids Res.*
379 2016;44(D1):D710–16.
- 380 32. Firestein S. How the olfactory system makes sense of scents. *Nature.*
381 2001;413(6852):211-8.
- 382 33. Ronquist F, Huelsenbeck JP. MrBayes 3: Bayesian phylogenetic inference
383 under mixed models. *Bioinformatics.* 2003;19(12):1572-4.
- 384 34. Guindon S, Dufayard JF, Lefort V, et al. New algorithms and methods to
385 estimate maximum-likelihood phylogenies: assessing the performance of
386 PhyML 3.0. *Syst Biol.* 2010;59(3):307-21.
- 387 35. Stamatakis A, Hoover P, Rougemont J. A rapid bootstrap algorithm for the
388 RAxML Web servers. *Syst Biol.* 2008;57(5):758-71.
- 389
390
391

Table 1 Genome sequencing information.

Insert size (bp)	Read length (bp)	Raw data		Clean data	
		Total bases (Gb)	Sequencing depth (x)	Total bases (Gb)	Sequencing depth (x)
230	125	135.76	46.02	125.96	42.70
500	125	102.51	34.75	88.52	30.01
2,000	125	59.0	20.00	50.16	17.00
5,000	125	51.57	17.48	46.39	15.73
10,000	125	28.16	9.55	24.67	8.36
15,000	125	30.34	10.28	28.14	9.54
Total		407.34	138.08	363.84	123.34

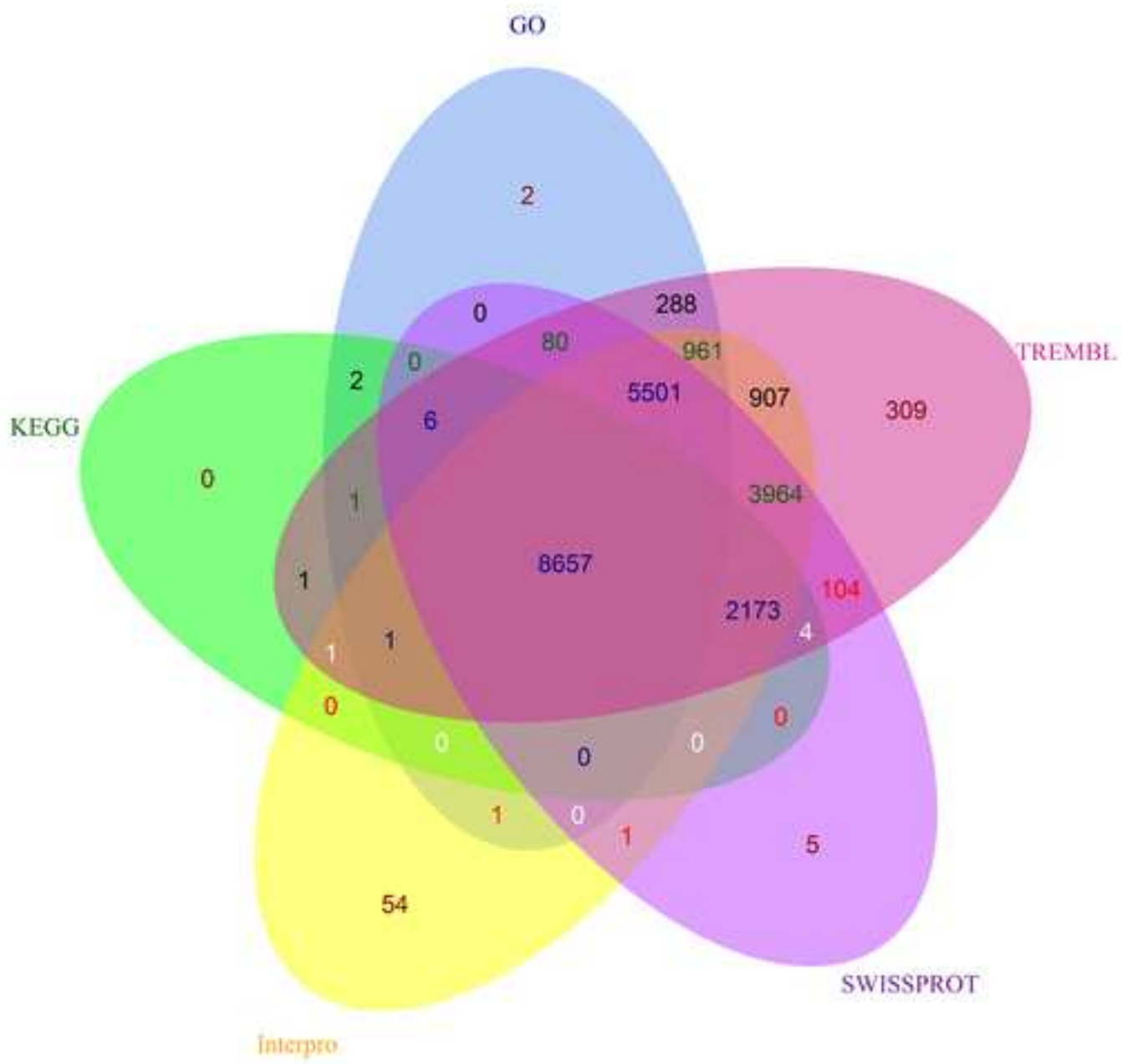
Note: Genome size is 2.95Gb.

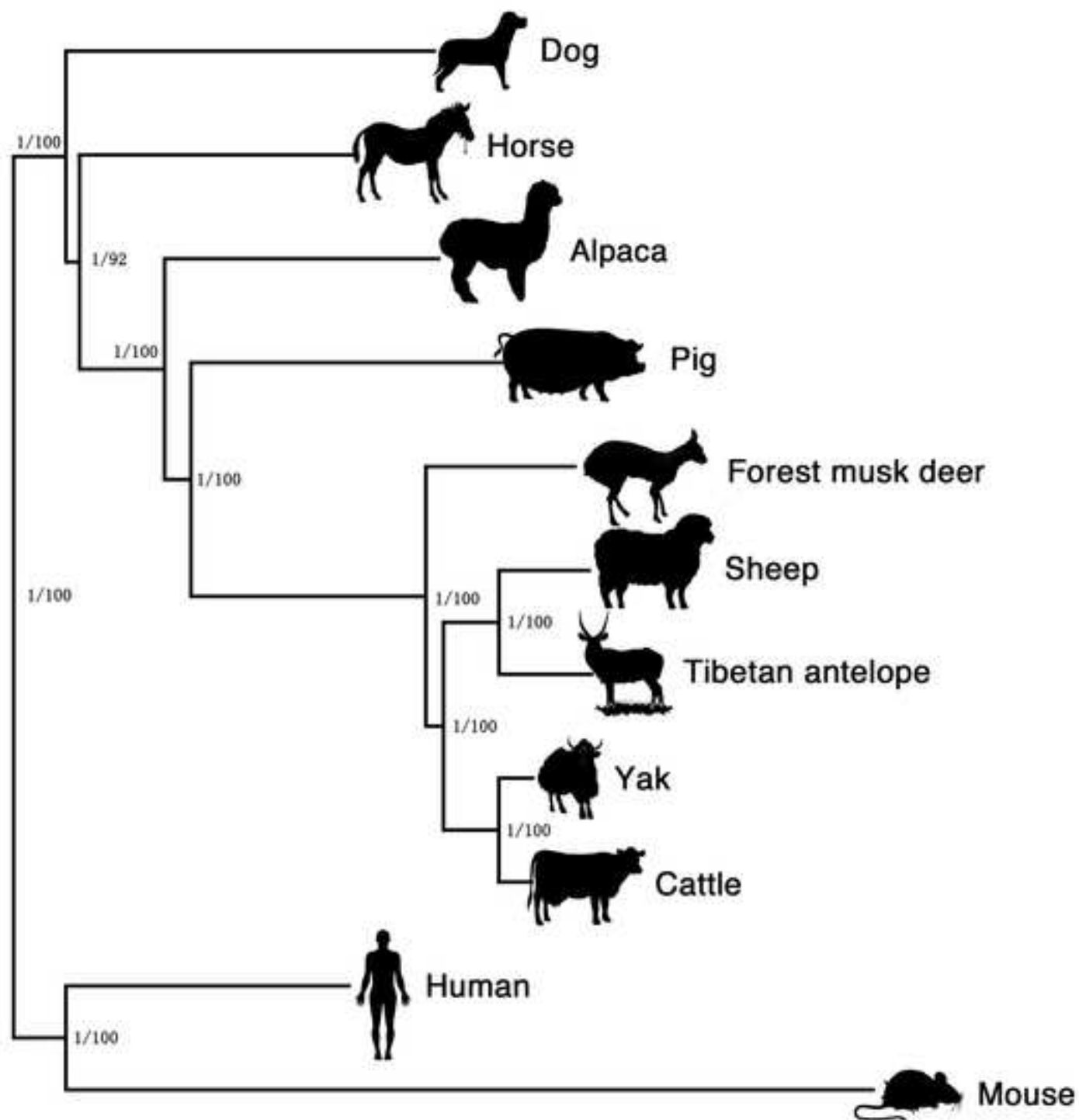
Table 2 Statistics of the final assembly of forest musk deer genome.

Genome assembly	Numbers
Contig N50 (kb)	22.6
Scaffold N50 (Mb)	2.85
Longest Scaffold(Mb)	18.69
Scaffold Number	79,206
GC content	40%
Total length (Gb)	2.72

Table 3 Functional annotation statistics of forest musk deer genome by various methods

	Database	Number	Percent (%)
Total		24,352	100.00
	Swissprot	18,771	77.08
	TrEMBL	22,696	93.20
Annotated	KEGG	10,846	44.54
	Interpro	22,221	91.12
	GO (blast2go)	15,736	64.62
	GO (Interproscan)	14,815	60.84
Un-annotated		1,329	5.77








Click here to access/download
Supplementary Material
MD_R1_suppNote.docx





Click here to access/download
Supplementary Material
Table S1_BUSCO.docx





Click here to access/download
Supplementary Material
Table S2_new.docx





Click here to access/download
Supplementary Material
Table S3_repeat.docx





Click here to access/download
Supplementary Material
Table S4_SSR.docx







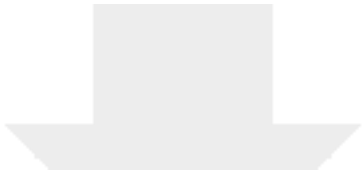
Click here to access/download
Supplementary Material
Figure S1_kmer.tif



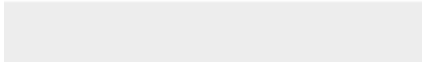


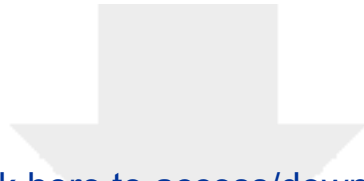
Click here to access/download
Supplementary Material
Figure S2_wego.pdf





Click here to access/download
Supplementary Material
Figure S3_TE_type.tif





Click here to access/download
Supplementary Material
Figure S4_OrthologousGene.tif



Dear Editor,

Please find attached our manuscript entitled “The draft genome sequence of forest musk deer (*Moschus berezovskii*)” that we would like considered for publication in *GigaScience*. We have carefully revised the manuscript in accordance with reviewers’ comments as detailed in the attached material. Although there were many comments requiring changes to the manuscript, we feel the comments and needed actions were straightforward and have thus been able to accommodate nearly every concern by changes to the text. We hope you will now find the manuscript acceptable for publication. We declare that the manuscript is not published, in press, or simultaneously submitted elsewhere. All authors have read and approved the manuscript, and declared that there were no competing interests.

Please address all correspondence concerning this manuscript to Dr Zhenxin Fan.

Sincerely,

Dr. Zhenxin Fan

College of Life Sciences, Sichuan University, Cheng Du, Sichuan 610065, China

E-mail: zxfan@scu.edu.cn