

GigaScience

Draft genome of the gayal, *Bos frontalis*

--Manuscript Draft--

Manuscript Number:	GIGA-D-17-00116R1	
Full Title:	Draft genome of the gayal, <i>Bos frontalis</i>	
Article Type:	Data Note	
Funding Information:	Chinese 973 program (2013CB835200, 2013CB835204)	Dr. Dong-Dong Wu
	Strategic Priority Research Program of the Chinese Academy of Sciences (XDB13020600)	Dr. Dong-Dong Wu
Abstract:	<p>Background: Gayal (<i>Bos frontalis</i>), also known as mithan or mithun, is a large endangered semi-domesticated bovine that has a limited geographical distribution in the hill-forests of China, Northeast India, Bangladesh, Myanmar, and Bhutan. Many questions about the gayal such as its origin, population history as well as genetic basis of local adaptation remain largely unresolved. De novo sequencing and assembly of the whole gayal genome provides an opportunity to address these issues.</p> <p>Findings: We report a high-depth sequencing, de novo assembly, and annotation of a female Chinese gayal genome. Based on the Illumina genomic sequencing platform, we have generated 350.38Gb raw data from 16 different insert-size libraries. A total of 276.86Gb clean data is retained after quality control. The assembled genome is about 2.85Gb with scaffold and contig N50 sizes of 2.74Mb and 14.41kb, respectively. Repetitive elements account for 48.13% of the genome. Gene annotation has yielded 26,667 protein-coding genes, of which 97.18% have been functionally annotated. BUSCO assessment shows that our assembly captures 93% (3,183 of 4,104) of the core eukaryotic genes, and 83.1% of vertebrate universal single-copy orthologs.</p> <p>Conclusions: We provide the first comprehensive de novo genome of the gayal. This genetic resource is integral for investigating the origin of the gayal and performing comparative genomic studies to improve understanding of the speciation and divergence of Bovine species. The assembled genome could be used as reference in future population genetic studies of gayal.</p>	
Corresponding Author:	Dong-Dong Wu CHINA	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:		
Corresponding Author's Secondary Institution:		
First Author:	Dong-Dong Wu	
First Author Secondary Information:		
Order of Authors:	Dong-Dong Wu	
	Ming-Shan Wang	
	Yan Zeng	
	Xiao Wang	
	Wen-Hui Nie	
	Jin-Huan Wang	

	Wei-Ting Su
	Newton O. Otecko
	Zi-Jun Xiong
	Sheng Wang
	Kai-Xing Qu
	Shou-Qing Yan
	Min-Min Yang
	Wen Wang
	Yang Dong
	Ya-Ping Zhang
Order of Authors Secondary Information:	
Response to Reviewers:	<p>Response to editor: Your manuscript "Draft genome of the gayal, <i>Bos frontalis</i>" (GIGA-D-17-00116) has been assessed by our reviewers. Although it is of interest, we are unable to consider it for publication in its current form. The reviewers have raised a number of points which we believe would improve the manuscript and may allow a revised version to be published in GigaScience. Their reports, together with any other comments, are below. In particular, the reviewers point out that the previous literature in this field must be referenced more completely and accurately. Given the previous, published work on genome sequencing of this species, you should also explain better what the novel contribution of your study is.</p> <p>Reply: Thanks the editor for handling our manuscript and comments. We found the comments and suggestions very helpful. We have revise the manuscripts carefully, improving literature review and referencing, as wells as clarifying sample origin and novelty of our research. We believe the revisions have greatly improved our manuscript for publication in your reputable journal.</p> <p>Please pay particular attention to point 7) of reviewer 2 regarding the inferred genome size - this may need careful re-assessment.</p> <p>Reply: Thanks for the comment. In the previous estimation, we used raw sequencing reads (without filtration) to infer the K-mer frequency and genome size. We have corrected this mistake and re-assessed the genome size using only the clean reads that passed quality filtration in the genome assembly. The newly estimated genome size is 3.15Gb, still slightly larger than what we assembled (2.85Gb). We have also illustrated and discussed in our responses the discrepancies that commonly occur between K-mer estimated and assembled genome sizes. Please full details in the response to point 4) of reviewer 1 and point 7) of reviewer 2.</p> <p>Please also provide more details regarding the origin of the sample, and address all other points of the reviewers.</p> <p>Reply: We have provided more details on sample origin. The gayal used in this study originated from Dulong, a city in Yunnan province, China. It is currently reared in Yunnan Academy of Grassland and Animal Science for breeding and research purposes. Karyotype examination showed it has 2n=58 chromosomes (see figure2). We have addressed all the points by the reviewers in the one by one response below.</p> <p>Response to reviewers: Reviewer #1: 1. Average exon in text is 3.27 where as in corresponding table it is 7.19</p> <p>Reply: We are very sorry for the mistake. We predicted genes using both homolog and de novo based methods. Both genes set were subsequently merged using glean to produce the final gene set, in which average exons per gene is 3.27. In the homolog method, using <i>Bos taurus</i> as closed species to search again gayal genome, we</p>

predicted 19,666 protein coding genes with average exons of 7.19 per gene. We have made appropriate revisions for consistency and clarity.

2. Reference for buffalo assembly is missing from references

Reply:

We are sorry for this oversight; we have added the reference accordingly.

3. References need to be rechecked as per text

Reply:

Thanks to the reviewer for the comment. We have carefully revised the references one by one.

4. There is a need to re-look into figure of 3.7gb as the genome size of Mithun

Reply:

Thanks for the comment. In the previous estimation, we used raw sequencing reads (without filtration) to infer the K-mer frequency and genome size. We have corrected this mistake and re-assessed the genome size using only the clean reads that passed quality filtration in the genome assembly. The newly estimated genome size is 3.15Gb, still slightly larger than what we assembled (2.85Gb). However, minimal discrepancies between K-mer estimated and assembled genome sizes is a common occurrence in NGS studies (Yim et al. Nat Genet. 2014;46(1):88-92; Wang et al. Gigascience. 2017;doi: 10.1093/gigascience/gix016; Fan et al. Nat Commun. 2013;4:1426; Gao et al. Gigascience. 2017; doi:10.1093/gigascience/gix041). We think low sequencing bases likely lead to over estimation of genome size. In addition, as demonstrated by the previous gayal sequencing (Mei et al.2016) and our current work, there is high heterozygosity in the gayal genome, which also likely influence its genome size estimation.

Reviewer #2:

This is a well-written account of a whole-genome sequence of the gayal, a most interesting bovine species. However, it should become clear what is the novelty of the results relative to an earlier report on a WGS of the same species. Furthermore, more details about the sample origin should be given, while referencing to the literature about the gayal is superficial and even incorrect. We recommend a major revision.

Reply:

We thank the reviewer very much for the constructive analysis and comments on our manuscript. We have followed the suggestions of the reviewer to revise our manuscript, particularly discussing the novelty of the results relative previous research on gayal and other bovine relatives, explaining sample origin, as well as revising the literature review and references. Please see below a detailed point by point response to the detailed comments.

Detailed comments

1. As cited, Mei et al. (2016) already published a gayal WGS, so a separate publication on another sequence should be justified, for instance because of a better coverage, contig and scaffold statistics and gene coverage.

Reply:

Thanks the reviewer for the comment. As stated by the reviewer, last year, Mei et al. published a study in which they re-sequenced gayal WGS. They generated 36.3Gb genome sequence data with an average sequencing depth of 13.06X after mapping the sequencing reads to cattle reference genome. Their analysis was therefore based on SNPs obtained by mapping gayal genome to cow reference genome. They further constructed phylogenetic trees using a subset of only 20 randomly selected single ortholog copy genes in *Bos taurus*, *Bos mutus* (wild yak) and *Bubalus bubalis* genomes, placing gayal off *B. mutus* and *B. taurus*. While we appreciate the importance of their work and other preceding partial genome research on gayal, we also take note that they used a resequencing approach for species that does not have a reference genome, forcing them to map the gayal sequencing reads to a cattle genome. Their study, as well as our own analysis, shows that gayal has a high heterozygosity and is far divergent from cattle. Hence, using cattle reference when mapping gayal sequencing reads is definitely likely to produce biases during alignment and SNP calling procedures. In addition, they did not determine/report the karyotype of the gayal they used. This is an important matter for ongoing research on gayal as

gayal hybrids are common in China. In our study, we have tried to take care of these limitations. We used a female gayal with 58 chromosomes to perform high coverage whole genome sequencing (350.38Gb raw data) with libraries constructed based on different insert sizes, and then performed de novo assembly. Besides the detailed analysis and description of the genome properties, we also state the karyotype of the gayal used and its phylogenetic relationship with other bovines (validated by complete mtDNA gayal sequences generated by Sanger sequencing method). Overall, our study represents the pioneer de novo assembly of the gayal whole genome, and Sanger sequencing of its complete mtDNA. Our study therefore presents a suitable reference genome for future studies on gayal, plus other important resources and insights that will facilitate research on gayal and other bovine species.

We have concisely included these descriptions in the revised manuscript.

2. The geographic origin sample of the sample should be specified. Chinese gayals, or Dulong cattle, are known to harbor zebu or taurine mtDNA (Gou et al. 2010, J.Anim.Breeding Genet. 127, 154-160; Mei et al. 2016) and may very well differ from individuals with an Indian origin.

Reply:

Thanks the reviewer for the comment. We have provided more details on sample origin. The gayal used in this study originated from Dulong, a city in Yunnan province, China. It is currently reared in Yunnan Academy of Grassland and Animal Science for breeding and research purposes. As suggested, we have explained the sample origin more clearly and cited these references appropriately in our revised manuscript.

3. For this reason the mtDNA sequence should be retrieved and compared to the several available gayal mtDNA sequences published previously.

Reply:

As suggested, we searched NCBI-Nucleotide database for published mtDNA sequences gayal. Unfortunately, there is no complete mtDNA assembly available for gayal, except partial mtDNA sequences like D-loop, cytb, and 16s. Considering the lower ability of NGS to accurately recover duplicated sequences that characterize regions like the D-loop in mtDNA, we sequenced complete mtDNA from the gayal in our study using Sanger method. We then downloaded sequences of mtDNA for gayal and other Bovine species, and constructed phylogenetic trees. Below are trees constructed using maximum likelihood method based on complete mtDNA (see figure 5) and cytb (see figure S4) sequences. We observed that the gayal in our study clustered with gaur and gayal from Dulong, Myanmar, Bhutan, and Manipuri. We have submitted the new complete mtDNA sequence to the Genbank and added this analysis in our revised manuscript.

4. Thai and Malaysian gaur have indeed a $2n=56$ karyotype, but Indian gaur, which occurs in the geographic area overlapping with the range of the gayals, has $2n=58$ (Winter et al., 1984, Res Vet Sci 36: 276-283; Gallagher et al., 1992, J Hered 83:287-298; Mastro Monaco et al., 2004, Chromosome Res. 2:725-31).

Reply:

We thanks the reviewer for this comment. We agree with the reviewer that determining and reporting the karyotype of gayal is important due to these cryptic variations. Besides reporting the karyotype of the gayal in our study, we have revised our manuscript to reflect the insights offered by the reviewer plus the appropriate citations.

5. The cited references (5,14) do not show that gaur x gayal male offspring are sterile. Although I could not find literature about the outcome of this hybrid cross, it is generally assumed that gayal is the domestic form of the gayal, also because they have similar mtDNA and Y-chromosomal DNA sequences (Hassanin et al. 2012, C.R.Biologies 335:32-50; Nijman et al. 2008, Cladistics 24:723-726).

Reply:

We are sorry for the oversight. We have revised this description to maintain only the details that have a solid literature backing. Thanks for the comment.

6. The URL reference [19] of the academic thesis describing the American bison WGS is still inaccessible. I guess that this WGS has been downloaded from Genbank, which should be made clear.

Reply:

	<p>Thanks to the reviewer for the comment. We download the sequence from Genbank and have revised the citation appropriately.</p> <p>7. The inferred genome size for the gayal of 3.7 Gbp, larger than the genome of any related mammalian species, is not believable and not consistent with the gene coverage. Reply: We thank the reviewer for this important observation. In the previous estimation, we used raw sequencing reads (without filtration) to infer the K-mer frequency and genome size. We have corrected this mistake and re-assessed the genome size using only the clean reads that passed quality filtration in the genome assembly. The newly estimated genome size is 3.15Gb, still slightly larger than what we assembled (2.85Gb). However, minimal discrepancies between K-mer estimated and assembled genome sizes is a common occurrence in NGS studies (Yim et al. Nat Genet. 2014;46(1):88-92; Wang et al. Gigascience. 2017;doi: 10.1093/gigascience/gix016; Fan et al. Nat Commun. 2013;4:1426; Gao et al. Gigascience. 2017; doi:10.1093/gigascience/gix041). We think that low quality sequencing bases likely lead to over estimation of genome size. In addition, as demonstrated by the previous gayal sequencing (Mei et al.2016) and our current work, there is high heterozygosity in the gayal genome, which also likely influence its genome size estimation.</p> <p>8. It may be interesting to compare the recovered DNA repeats with those from the bovine WGS. Reply: Thanks to the reviewer for the comment. It is an interesting topic to compare the repeats in different bovine species. However, whole genome solely based on NGS has low efficiency to assemble repeat sequences (Wang et al.2016. Nature Genetics 48(9): 972-3). In addition, many of these bovine genomes are generated without uniform sequencing platform and assembly strategies. Further, the repeats predictions do not follow a harmonized pipeline, hence remain just draft genomes. It is difficult to distinguish the lose or increase of repeats in one species to be attributable to evolution or from technique/sequencing effects. The main reach of the current study is providing a comprehensive genetic resources and a draft reference genome for gayal to facilitate future research. We believe that in future, when high quality genomes for the bovine species become available, it will be fascinating to retrieve and compare DNA repeats evolution among the bovine species.</p> <p>9. page 9 last line: vertebrata > vertebrate. Reply: Thanks, we have revised this accordingly.</p>
Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
<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
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>

Data Note

Draft genome of the gayal, *Bos frontalis*

Ming-Shan Wang^{1,2,#}, Yan Zeng^{1,2,#}, Xiao Wang^{1,2,#}, Wen-Hui Nie¹, Jin-Huan Wang¹, Wei-Ting Su¹, Newton O. Otecko^{1,2}, Zi-Jun Xiong^{1,4}, Sheng Wang⁵, Kai-Xing Qu⁶, Shou-Qing Yan⁷, Min-Min Yang^{1,2}, Wen Wang^{1,2}, Yang Dong^{8,9,*}, Dong-Dong Wu^{1,2*}, and Ya-Ping Zhang^{1,2,3,*}

1. State Key Laboratory of Genetic Resources and Evolution, Yunnan Laboratory of Molecular Biology of Domestic Animals, Kunming Institute of Zoology, Chinese Academy of Sciences, Kunming 650223, China
2. Kunming College of Life Science, University of Chinese Academy of Sciences, Kunming 650204, China
3. Laboratory for Conservation and Utilization of Bio-resource, Yunnan University, Kunming 650091, China
4. China National GeneBank, BGI-Shenzhen, Shenzhen 518083, China
5. Laboratory of Animal Genetics, Breeding and Reproduction, Ministry of Agriculture of China, National Engineering Laboratory for Animal Breeding, College of Animal Science and Technology, China Agricultural University, Beijing 100193, China
6. Yunnan Academy of Grassland and Animal Science, Kunming 650212, China
7. College of Animal Science, Jilin University, Changchun 130062, China
8. Yunnan Agricultural University, Kunming 650100, China
9. Faculty of Life Science and Technology, Kunming University of Science and Technology, Kunming 650500, China

Contribute equally to this work

* Corresponding author:

Dr. Yang Dong, email: loyalyang@163.com

Dr. Dong-Dong Wu, email: wudongdong@mail.kiz.ac.cn

Dr. Ya-Ping Zhang, email: zhangyp@mail.kiz.ac.cn

ORCID details:

Newton O. Otecko: 0000-0002-9149-4776

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

1 **Abstract**

2 **Background:**

3
4 Gayal (*Bos frontalis*), also known as mithan or mithun, is a large endangered
5 semi-domesticated bovine that has a limited geographical distribution in the
6 hill-forests of China, Northeast India, Bangladesh, Myanmar, and Bhutan. Many
7 questions about the gayal such as its origin, population history as well as genetic basis
8 of local adaptation remain largely unresolved. *De novo* sequencing and assembly of
9 the whole gayal genome provides an opportunity to address these issues.

10 **Findings:**

11 We report a high-depth sequencing, *de novo* assembly, and annotation of a female
12 Chinese gayal genome. Based on the Illumina genomic sequencing platform, we have
13 generated 350.38Gb raw data from 16 different insert-size libraries. A total of
14 276.86Gb clean data is retained after quality control. The assembled genome is about
15 2.85Gb with scaffold and contig N50 sizes of 2.74Mb and 14.41kb, respectively.
16 Repetitive elements account for 48.13% of the genome. Gene annotation has yielded
17 26,667 protein-coding genes, of which 97.18% have been functionally annotated.
18 BUSCO assessment shows that our assembly captures 93% (3,183 of 4,104) of the
19 core eukaryotic genes, and 83.1% of vertebrate universal single-copy orthologs.

20 **Conclusions:**

21 We provide the first comprehensive *de novo* genome of the gayal. This genetic
22 resource is integral for investigating the origin of the gayal and performing
23 comparative genomic studies to improve understanding of the speciation and
24 divergence of Bovine species. The assembled genome could be used as reference in
25 future population genetic studies of gayal.

26 **Keywords:** *Bos frontalis*; Genome assembly; Annotation; Phylogeny
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

Data description

Background

The gayal is a large-sized endangered semi-domesticated bovine species belonging to the family Bovidae, tribe Bovini, group Bovina, genus *Bos* and species *Bos frontalis* (NCBI Taxon ID: 30520). It is also called the mithan or mithun. Its distribution spans eastern Bhutan through the Arunachal Pradesh in India to the Naga and Chin hills in the Arakan Yomarang region that defines the borders between India, Bangladesh, Myanmar, and China [1, 2]. The Gayal has unique characters and appearances compared to gaur, cattle, and other bovine species [3]. These features include a bony dorsal ridge on the shoulder and white stockings on all four legs (**Figure 1**). It has been previously held that gayal was domesticated from gaur and/or from a hybrid descendant from crossing domestic cattle (*B. indicus* or *B. taurus*) and wild gaur [2, 4, 5]. Karyotype analysis indicates that Indian gayal has a $2n=58$ karyotype, same as the local gaur ($2n=58$) [6, 7], but different from Chinese and Malaysian gaurs (*B. gaurus*, $2n=56$) as well as domesticated cattle (*B. indicus* and *B. taurus*, $2n=60$) [2, 6-10].

Phylogenetic analyses in multiple studies based on mtDNA or Y-chromosomal DNA place gayal in conflicting clustering positions with respect to cattle, zebu and wild gaur. For example, Chinese gayal, or Dulong cattle, are known to harbor zebu or taurine mtDNA footprints, suggesting hybrid origin [5, 11]; and more studies have shown a high mtDNA and Y-chromosomal DNA sequences similarity between gayal and guar [12-15]. One study has even placed the gayal as a distinct and separate species/sub-species [16]. In contrast, phylogenetic analyses based on SNPs from 20

1 randomly selected single copy gene orthologs of *B. taurus*, *B. mutus* (wild yak) and
2
3 *Bubalus bubalis* placed Chinese gayal off the *B. mutus* and *B. taurus* clade, indicating
4
5 that gayal is distinct from the modern domestic cattle, *B. taurus* [5]. These authors
6
7 further demonstrated from mtDNA analysis that the gayal is the most proximal to
8
9 domesticated cattle (*B. taurus* and *B. indicus*), suggesting that the gayal could be a
10
11 hybrid emanating from crossing of male wild gaur and female domestic cattle [5].
12
13
14 These differences illustrate the existence of unresolved uncertainties regarding the
15
16
17 origin of gayal.
18
19
20
21
22

23 Research has revealed a high genomic divergence among bovine species [17, 18].

24
25
26 Consequently, mapping of resequencing data from one bovine species onto the
27
28
29 reference genome of different species (for instance, gayal versus cattle) creates
30
31
32
33 avenues for biases and/or errors in sequence alignment and SNP calling procedures.
34
35

36 This challenge extends to species of great research interest like gayal, which so far
37
38
39 have no *de novo* assembled reference genome. For instance, Mei et al. recently
40
41
42 reported a whole genome sequencing (resequencing) of Chinese gayal [5]. In their
43
44
45 analysis, they retrieved variants based on mapping gayal sequencing reads (13.06X)
46
47
48 to the cattle reference genome. Importantly, hybrid gayals are hard to distinguish
49
50
51 only through morphological characterization, yet Mei *et al.* did not examine the
52
53
54 karyotype of the gayal they resequenced. In contrast to the gayal, *de novo* genome
55
56
57 assembly has been accomplished for related species like cattle (*Bos taurus*) [19], yak
58
59
60
61
62
63
64
65

1
2 (*Bos grunniens*) [17], wisent (*Bison bonasus*) [20], North American bison (*Bison*
3
4 *bison*) [21], zebu (*Bos indicus*) [22], and water buffalo (*Bubalus bubalis*) [23]. This
5
6
7 represents a critical resource towards mitigating the challenges inherent in
8
9
10 resequencing approaches, and provides great opportunities to refine the evolutionary
11
12 history of bovine species. In this study, we for the first time report the draft genome
13
14 assembly of the gayal with a high sequencing depth generated on the Illumina genome
15
16
17 sequencing platform. This valuable resource is an important input to the research of
18
19
20 the origin and evolution of this species that has been classified as an endangered by
21
22
23 the IUCN.
24
25
26
27
28
29

30 **Sample collection and sequencing**

31
32
33
34 The gayal (NCBI taxonomy ID: 30520) used for genome sequencing came from a
35
36
37 Dulong in Yunnan province, China (**Figure 1**). It was kept at Yunnan Academy of
38
39
40 Grassland and Animal Science for breeding and research purposes. Karyotype
41
42
43 examination showed that it has $2n=58$ chromosomes (**Figure 2**). We extracted total
44
45
46 genomic DNA from skin fibroblast cell lines of the gayal using Qiagen Blood and
47
48
49 Tissue Kit (Qiagen, Valencia, CA) according to the manufacturer's instructions. The
50
51
52 cells are maintained at the Cell Bank of Kunming Institute of Zoology (specimen ID:
53
54
55 KCB201042). A total of 17 paired-end genomic sequence libraries were constructed
56
57
58 with a gradient insert size ranging from 180bp to 20kb, and sequencing was carried
59
60
61 out on the Illumina HiSeq 2000 platform according to the manufacturer's instructions.
62
63
64
65

1 For short insert size libraries (180bp, 250bp, 450 bp, and 600bp), sequencing was
2
3 performed at the Central Laboratory of Kunming Institute of Zoology with read
4
5 lengths of 100bp. Sequencing of long insert size libraries (800 bp, 2, 5, 10 and 20 kb)
6
7 was conducted at BGI-Shenzhen with read lengths of 49bp, except for the 800bp
8
9 insert size library which were sequenced with a read length of 85bp. A total of
10
11 350.38Gb raw sequence data has been generated in our study (**Additional file 1:**
12
13 **Table S1**). Before assembly, we performed strict quality control by removing poor
14
15 quality reads and/or bases using scripts from SOAPec (version 2.02) [24]. Reads were
16
17 shortened by 2bp at both head and tail. We dropped any read plus their corresponding
18
19 paired-end if it contained more than 30 low-quality bases or more than 5% unknown
20
21 base (usually denoted by N). Reads with duplications and adapters were also removed.
22
23 We corrected for sequencing errors using the k-mer (13 used in this study) frequency
24
25 method in SOAPec (version 2.02) [24]. After filtering and correction, we retained
26
27 276.86Gb high-quality sequences for genome assembly (**Additional file 1: Table S2**).

40 ***De novo* assembly of gayal genome**

41
42 In order to have a basic knowledge about the genome size and attributes of the gayal
43
44 genome, we performed a 17-mer analysis using clean and high quality sequences from
45
46 180 and 450bp insert size libraries. We extracted the 17-mer sequences using sliding
47
48 windows with a size of 17bp, and calculated the frequency of each 17-mer. A clear
49
50 peak at 25X with two upward convex signals besides it is evident, suggesting high
51
52 heterozygosity. The genome size for gayal is estimated to be 3.15Gb (**Additional file**
53
54 **1: Table S3; Figure 3**).

1 We then performed *de novo* assembly of gayal genome using Platanus (version
2
3 2.0) (Platanus, RRID:SCR_015531) [25] in three steps: contig construction,
4
5 scaffolding, and gap filling. To construct contigs based on short insert size libraries
6
7 (180, 250, 450, 600 and 800bp), we used Platanus (version 2.0) [25], which includes a
8
9 series of procedures such as constructing de Bruijn graphs, clipping tips, merging
10
11 bubbles, and removing low coverage links. In the scaffolding step, reads from both
12
13 small and large insert libraries were mapped to contig sequences to construct scaffolds
14
15 using distance information from read pairs. An additional local assembly of reads,
16
17 with one end of a read pair uniquely aligned to a contig and the other end located
18
19 within the gap, was performed using GapCloser (version 1.12) (GapCloser ,
20
21 RRID:SCR_015026) [24]. These processes yielded a final draft gayal genome
22
23 assembly with a total length of 2.85Gb, contig N50 of 14.4 kb, and scaffold N50 of
24
25 2.74Mb (**Table 1**). The assembled genome size is similar to that reported for cattle
26
27 [26] and yak [17]. To assess the completeness of the assembled gayal genome, we
28
29 performed BUSCO analysis (BUSCO , RRID:SCR_015008) [27] by searching
30
31 against the arthropod universal benchmarking single-copy orthologs (BUSCOs,
32
33 version 2.0). Overall, 85.2% and 7.8 % of the 4,104 expected vertebrate genes are
34
35 identified in the assembled genome as complete and partial, respectively.
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

Approximately 291 genes could be considered missing in our assembly. Of the
expected complete vertebrate genes, 3434 and 60 are identified as single copy and
duplicated BUSCOs, respectively (**Table 2**). Our newly assembled gayal genome has
a slightly lower completeness rate compared to genomes of yak [17], wisent [20],

1 bison [21], zebu [22], and buffalo [23] (**Table 2**).

6 **Annotation of genomic repeat sequences in gayal genome**

7
8
9 To search for the repeated sequences in gayal genome, including tandem repeats,
10 interspersed repeats, and transposable elements (TE) (e.g. LINE, SINE, LTR, DNA
11 transposons), we leveraged both *de novo* and homolog-based methods as used in
12 previous publications [28, 29]. For the homolog-based methods, we used
13 RepeatMasker (RepeatMasker , RRID:SCR_012954) and RepeatProteinMask
14 (<http://repeatmasker.org/>) to search against the known Repbase TE library
15 (RepBase21.01) [30] and TE protein database, respectively. In the *de novo* method,
16 Piler [31] and RepeatModeler (RepeatModeler , RRID:SCR_015027)
17 (<http://www.repeatmasker.org/>) are used to generate a *de novo* gayal repeat library,
18 which is subsequently used in Repeat-Masker to annotate repeats. TRF [32] is then
19 employed to predict tandem repeats. The combined results show that a total of 1.37Gb
20 non-redundant repetitive sequences are identified in the gayal genome, which account
21 for 48.13% of the whole genome. The most predominant repeat is the long
22 interspersed nuclear elements (LINEs), which account for 40.43% (1.15Gb in total) of
23 the genome (**Table 3; Additional file 1: Table S4, Figure S1, Figure S2**).

53 **Gayal genome gene structure prediction**

54
55 For gene structure prediction, we combined both *de novo* and homolog-based
56 approaches to predict protein-coding genes in the gayal genome. In homolog-based
57
58
59

1 method, gene sets from *Bos taurus* [19], *Canis familiaris* [33], *Homo sapiens*
2
3 (ENSEMBL 80), *Sus scrofa* [34], *Rattus norvegicus* (ENSEMBL 80), and *Ovis aries*
4
5 [35] were used as queries to search against gayal genome (**Additional file 1: Table**
6
7 **S5**). For the *de novo* based method, AUGUSTUS (Augustus: Gene Prediction ,
8
9 RRID:SCR_008417) [36], Genescan (GENSCAN , RRID:SCR_012902) [37], and
10
11 GlimmerHMM (GlimmerHMM , RRID:SCR_002654) [38] were used as engines to
12
13 predict gene models. We then merged the gene prediction results derived from both
14
15 methods using GLEAN [39] to generate a consensus gene set. In total, we have
16
17 identified 26,667 protein coding genes with a mean of 3.27 exons per gene (**Table 4;**
18
19 **Additional file 1: Figure S3**). The lengths of genes, CDS, introns, and exons in gayal
20
21 are comparable to those of the genomes used for homolog-based predictions
22
23 (**Additional file 1: Figure S3**). In addition, we predicted non-coding RNA genes in
24
25 the gayal genome. We used blast to search rRNA against Human rRNA database, and
26
27 tRNAscan-SE (tRNAscan-SE , RRID:SCR_010835) [40] to search tRNA in the
28
29 genome sequences. We also used blast to search miRNA and snRNA via Rfam
30
31 database (Rfam , RRID:SCR_007891)(release 11.0) [41]. We reveal a total of 2,357
32
33 ribosomal RNA (rRNA), 29,821 transfer RNA (tRNA), 16,305 microRNAs (miRNA),
34
35 and 1,380 snRNA genes in the gayal genome (**Additional file 1: Table S5**).

52 **Functional annotation of protein-coding genes**

53
54 Gene functional annotation refers to searching functional motifs, domains, and
55
56 possible biological process by aligning translated gene coding sequences to known
57
58
59
60
61
62
63
64
65

1 databases such as SwissProt and TrEMBL [42], NT database (from NCBI), Gene
2
3 Ontology (GO) (GO , RRID:SCR_002811), and Kyoto Encyclopaedia of Genes and
4
5 Genomes (KEGG , RRID:SCR_012773)(KEGG) [43]. We have annotated all the
6
7 protein coding genes identified in this study to retrieve functional terms according to
8
9 InterPro, KEGG, and GO terms. Overall, 81.74% (21,798), 54.56% (14,550), and
10
11 66.39% (17,704) genes show enrichment in InterPro, KEGG, and GO respectively. In
12
13 total, 25,916 protein-coding genes (97.18%) were successfully annotated for
14
15 conserved functional motifs and functional terms (**Additional file 1: Table S6**).
16
17
18
19
20
21
22
23
24

25 **Phylogenetic analysis and divergence time estimation**

26
27 To investigate the phylogenic position of gayal, we retrieved nucleotide and protein
28
29 data for cattle (*Bos taurus*) [19], yak (*Bos grunniens*) [17], wisent (*Bison bonasus*)
30
31 [20], bison (*Bison bison*) [21], zebu (*Bos indicus*) [22], and buffalo (*Bubalus bubalis*)
32
33 [23] from the NCBI database. Gene ortholog relationships of gayal and other bovine
34
35 species were identified by reciprocal blast searching with an e-value of 1e-7. Genes
36
37 with alternative splicing variants are represented by the longest transcript. Multiple
38
39 sequence alignment of the genes within one copy gene sets were performed using
40
41 MUSCLE program (MUSCLE , RRID:SCR_011812) [44]. Aligned sequences were
42
43 trimmed to remove potentially unreliably aligned regions and gaps using Gblocks [45].
44
45 Alignments with lengths shorter than 100bp were also discarded. Four-fold
46
47 degenerate sites were extracted and concatenated into a supergene. Modeltest [46]
48
49 was used to select the best substitution model. MrBayes (MrBayes ,
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 RRID:SCR_012067) [47] and RaxML (RAxML , RRID:SCR_006086) [48] software
2
3 were used to reconstruct the evolutionary relationships between species, and MEGA5
4
5 [49] used to view the tree. From these analyses, gayal clusters with the common
6
7 ancestor of cattle and zebu (**Figure 4**).

11 Additionally, we sequenced the complete mitochondrial DNA (mtDNA, the first
12
13 complete mtDNA of the gayal submitted to GenBank: MF614103) using Sanger
14
15 sequencing method, due to the fact that next generation sequencing methods have
16
17 lower ability and accuracy in recovering repeat sequences [28, 50], particularly in
18
19 regions with rich GC content like the D-loop. We then downloaded mtDNA sequences
20
21 of gayal and other bovine species from GenBank for phylogenic analysis. As shown in
22
23 **Figure 5** and **Figure S4**, the gayal we sequenced clusters with gaur (**Figure 5**,
24
25 **Additional file 1: Figure S4**). Our results from both whole genome and mtDNA data
26
27 differ from the conclusion made by Mei *et al.* who mapped gayal genome
28
29 resequencing data to a bovine reference [5]. Furthermore, the MCMCTREE program,
30
31 implemented using the PAML (PAML , RRID:SCR_014932)[51] package, was used
32
33 to estimate divergence times. The JC69 model and correlated molecular clock rates
34
35 (clock=3) were used in the calculation. Calibration time for the common ancestor of
36
37 buffalo and cattle obtained from the TimeTree database (<http://www.timetree.org/>)
38
39 was used to calibrate the divergence time. This analysis estimated the divergence time
40
41 of gayal from cattle and zebu at approximately 5.1 million years ago (**Figure 6**).

58 In conclusion, we have constructed a *de novo* assembly of the gayal genome and
59
60
61
62
63
64
65

1 describe its genetic attributes. To our knowledge, this is the first *de novo* assembled
2
3 genome for this species. We also demonstrate that together with the genomes of other
4
5 bovine species, the new gayal genome supports investigations concerning the origin,
6
7 evolutionary history, and local adaptation of gayal. This resource is also important for
8
9 the future conservation of this endangered species. In addition, the *de novo* gayal
10
11 genome adds to the list of available bovine genomes, and has advantages over
12
13 resequenced genomes in allowing accurate whole genome alignment, retrieving
14
15 constraint and/or rapidly evolved elements. It also strengthens the capacity to better
16
17 assess introgression, incomplete lineage sorting (ILS), and structural variation (SV)
18
19 among the bovine species, as well as inferring their effects on the species tree. The
20
21 assembled genome could be used as a reference in population genomic studies [52] of
22
23 the gayal. Furthermore, comprehensive comparative analyses of these genomes will
24
25 improve understanding of the formation and speciation of bovine species.
26
27
28
29
30
31
32
33
34
35
36
37
38

39 **Availability of supporting data**

40
41
42 The genome sequencing raw reads were deposited in the NCBI SRA database, project
43
44 ID: PRJNA387130. The assembly and annotation of the gayal genome are available in
45
46 the *GigaScience* GigaDB database[53]. The complete mtDNA for the gayal generated
47
48 by Sanger sequencing is also available in GenBank under the ID: MF614103. All
49
50 supplementary figures and tables are provided in Additional file 1.
51
52
53
54
55
56
57

58 **Competing interests**

59
60
61
62
63
64
65

1 The authors declare that they have no competing interests.
2
3
4
5

6 **Authors' contributions**

7

8
9 YPZ, DDW and MSW designed the study. WW and YD supervised the analyses.
10
11 WHN, WTS and JHW cultivated the cells. YZ and XW performed genome assembly
12
13 and annotation. MSW extracted genomic DNA and wrote manuscript with other
14
15 author's input. MSW and SQY sequenced the gayal complete mitochondrial DNA and
16
17 submitted to GenBank. SW, ZJX, KXQ, NOO, DY, DDW and YPZ revised the
18
19 manuscript. All authors read and approved the final manuscript.
20
21
22
23
24
25
26
27

28 **Acknowledgements**

29

30
31 This work was supported by the Strategic Priority Research Program of the Chinese A
32
33 cademy of Sciences, Grant No. XDB13020600, and the Chinese 973 program (2013C
34
35 B835200, 2013CB835204). DDW was supported by the Youth Innovation Promotion
36
37 Association, Chinese Academy of Sciences. N.O.O. is thankful for the support of the
38
39 CAS-TWAS President's Fellowship Program for Doctoral Candidates.
40
41
42
43
44
45
46
47

48 **References**

49

- 50 1. Payne WJA, Hodges J: **Tropical cattle: origins, breeds and breeding policies.**
51 1997:Blackwell Science.
52
- 53 2. Uzzaman MR, Bhuiyan MS, Edea Z, Kim KS: **Semi-domesticated and Irreplaceable**
54 **genetic resource gayal (*Bos frontalis*) needs effective genetic conservation in**
55 **Bangladesh: a review.** *Asian-Australas J Anim Sci* 2014, **27**:1368-1372.
56
57
- 58 3. Miao YW, Ha F, Gao HS, Yuan F, Li DL, Yuan YY: **Polymorphisms of inhibin α gene**
59
60
61
62
63
64
65

1 **exon 1 in buffalo (*Bubalus bubalis*), gayal (*Bos frontalis*) and yak (*Bos grunniens*).**
2 *Zool Res* 2012, **33**:402-408.
3

- 4 4. Payne WJA: **Cattle production in the tropics. Vol. 1. General introduction and breeds**
5 **and breeding.** London: Longman Group Ltd.; 1970.
6
7
8 5. Mei C, Wang H, Zhu W, Wang H, Cheng G, Qu K, Guang X, Li A, Zhao C, Yang W, et al:
9 **Whole-genome sequencing of the endangered bovine species gayal (*Bos frontalis*)**
10 **provides new insights into its genetic features.** *Sci Rep* 2016, **6**:19787.
11
12 6. Winter H, Mayr B, Schleger W, Dworak E, Krutzler J, Burger B: **Karyotyping, red**
13 **blood cell and haemoglobin typing of the mithun (*Bos frontalis*), its wild ancestor**
14 **and its hybrids.** *Res Vet Sci* 1984, **36**:276-283.
15
16
17 7. Gallagher DS, Jr., Womack JE: **Chromosome conservation in the Bovidae.** *J Hered*
18 1992, **83**:287-298.
19
20
21 8. Shan XN, Chen YF, Luo LH, Cao XM, Song JZ, Zeng YZ: **Comparative studies on the**
22 **chromosomes of five species of catties of the genus *Bos* in China.** *Zool Res* 1980,
23 **1**:75-81.
24
25 9. Adbullah MH, Idris I, Hilmi M: **Karyotype of malayan gaur (*Bos gaurus hubbacki*),**
26 **Sahiwal-Friesian cattle and gaur x cattle hybrid backcrosses.** *Pak J Biol Sci* 2009,
27 **12**:896-901.
28
29 10. Qu KX, He ZX, Nie WH, Zhang JC, Jin XD, Yang GR, Yuan XP, Huang BZ, Zhang YP,
30 Zhan LS: **Karyotype analysis of mithun (*Bos frontalis*) and mithun bull x Brahman**
31 **cow hybrids.** *Genet Mol Res* 2012, **11**:131-140.
32
33 11. Gou X, Wang Y, Yang S, Deng W, Mao H: **Genetic diversity and origin of gayal and**
34 **cattle in Yunnan revealed by mtDNA control region and SRY gene sequence**
35 **variation.** *J Anim Breed Genet* 2010, **127**:154-160.
36
37 12. Li SP, Chang H, Ma GL, Chen HY, Ji DJ, Geng RQ: **Molecular phylogeny of the gayal**
38 **inferred from the analysis of cytochrome b gene entire sequences.** *Yi Chuan* 2008,
39 **30**:65-70.
40
41 13. Nijman IJ, van Boxtel DCJ, van Cann LM, Marnoch Y, Cuppen E, Lenstra JA:
42 **Phylogeny of Y chromosomes from bovine species.** *Cladistics* 2008, **24**:723-726.
43
44 14. Dorji T, Mannen H, Namikawa T, Inamura T, Kawamoto Y: **Diversity and phylogeny of**
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 **mitochondrial DNA isolated from mithun *Bos frontalis* located in Bhutan.** *Anim*
2 *Genet* 2010, **41**:554-556.

- 3
4 15. Tanaka K, Takizawa T, Murakoshi H, Dorji T, Nyunt MM, Maeda Y, Yamamoto Y,
5 Namikawa T: **Molecular phylogeny and diversity of Myanmar and Bhutan mithun**
6 **based on mtDNA sequences.** *Anim Sci J* 2011, **82**:52-56.
7
8 16. Baig M, Mitra B, Qu KX, Peng MS, Ahmed I, Miao YW, Zan LS, Zhang YP:
9 **Mitochondrial DNA diversity and origin of *Bos frontalis*.** *Current Science* 2013,
10 **104**:115-120.
11
12 17. Qiu Q, Zhang G, Ma T, Qian W, Wang J, Ye Z, Cao C, Hu Q, Kim J, Larkin DM, et al:
13 **The yak genome and adaptation to life at high altitude.** *Nat Genet* 2012, **44**:946-949.
14
15 18. Porto-Neto LR, Sonstegard TS, Liu GE, Bickhart DM, Da Silva MV, Machado MA,
16 Utsunomiya YT, Garcia JF, Gondro C, Van Tassell CP: **Genomic divergence of zebu and**
17 **taurine cattle identified through high-density SNP genotyping.** *BMC Genomics* 2013,
18 **14**:876.
19
20 19. Zimin AV, Delcher AL, Florea L, Kelley DR, Schatz MC, Puiu D, Hanrahan F, Pertea G,
21 Van Tassell CP, Sonstegard TS, et al: **A whole-genome assembly of the domestic cow,**
22 ***Bos taurus*.** *Genome Biol* 2009, **10**:R42.
23
24 20. Wang K, Wang L, Lenstra JA, Jian J, Yang Y, Hu Q, Lai D, Qiu Q, Ma T, Du Z, et al: **The**
25 **genome sequence of the wisent (*Bison bonasus*).** *Gigascience*. 2017 Mar 10. doi:
26 10.1093/gigascience/gix016.
27
28 21. **American bison (*Bison bison bison*) genome assembly**
29 https://www.ncbi.nlm.nih.gov/assembly/GCF_000754665.1/.
30
31 22. Canavez FC, Luche DD, Stothard P, Leite KR, Sousa-Canavez JM, Plastow G, Meidanis J,
32 Souza MA, Feijao P, Moore SS, Camara-Lopes LH: **Genome sequence and assembly of**
33 ***Bos indicus*.** *J Hered* 2012, **103**:342-348.
34
35 23. **Water buffalo (*Bubalus bubalis*) genome assembly**
36 https://www.ncbi.nlm.nih.gov/assembly/GCA_000471725.1#/st.
37
38 24. Luo R, Liu B, Xie Y, Li Z, Huang W, Yuan J, He G, Chen Y, Pan Q, Liu Y, et al:
39 **SOAPdenovo2: an empirically improved memory-efficient short-read de novo**
40 **assembler.** *Gigascience* 2012, **1**:18.
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
25. Kajitani R, Toshimoto K, Noguchi H, Toyoda A, Ogura Y, Okuno M, Yabana M, Harada M, Nagayasu E, Maruyama H, et al: **Efficient de novo assembly of highly heterozygous genomes from whole-genome shotgun short reads.** *Genome Res* 2014, **24**:1384-1395.
 26. Bovine Genome S, Analysis C, Elisk CG, Tellam RL, Worley KC, Gibbs RA, Muzny DM, Weinstock GM, Adelson DL, Eichler EE, et al: **The genome sequence of taurine cattle: a window to ruminant biology and evolution.** *Science* 2009, **324**:522-528.
 27. Simao FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM: **BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs.** *Bioinformatics* 2015, **31**:3210-3212.
 28. Wang MS, Li Y, Peng MS, Zhong L, Wang ZJ, Li QY, Tu XL, Dong Y, Zhu CL, Wang L, et al: **Genomic analyses reveal potential independent adaptation to high altitude in Tibetan chickens.** *Mol Biol Evol* 2015, **32**:1880-1889.
 29. Xiong Z, Li F, Li Q, Zhou L, Gamble T, Zheng J, Kui L, Li C, Li S, Yang H, Zhang G: **Draft genome of the leopard gecko, *Eublepharis macularius*.** *Gigascience* 2016, **5**:47.
 30. Kapitonov VV, Jurka J: **A universal classification of eukaryotic transposable elements implemented in Repbase.** *Nat Rev Genet* 2008, **9**:411-412.
 31. Edgar RC, Myers EW: **PILER: identification and classification of genomic repeats.** *Bioinformatics* 2005, **21 Suppl 1**:152-158.
 32. Benson G: **Tandem repeats finder: a program to analyze DNA sequences.** *Nucleic Acids Res* 1999, **27**:573-580.
 33. Lindblad-Toh K, Wade CM, Mikkelsen TS, Karlsson EK, Jaffe DB, Kamal M, Clamp M, Chang JL, Kulbokas EJ, 3rd, Zody MC, et al: **Genome sequence, comparative analysis and haplotype structure of the domestic dog.** *Nature* 2005, **438**:803-819.
 34. Groenen MA, Archibald AL, Uenishi H, Tuggle CK, Takeuchi Y, Rothschild MF, Rogel-Gaillard C, Park C, Milan D, Megens HJ, et al: **Analyses of pig genomes provide insight into porcine demography and evolution.** *Nature* 2012, **491**:393-398.
 35. Jiang Y, Xie M, Chen W, Talbot R, Maddox JF, Faraut T, Wu C, Muzny DM, Li Y, Zhang W, et al: **The sheep genome illuminates biology of the rumen and lipid metabolism.** *Science* 2014, **344**:1168-1173.
 36. Stanke M, Keller O, Gunduz I, Hayes A, Waack S, Morgenstern B: **AUGUSTUS: ab**

1 **initio prediction of alternative transcripts.** *Nucleic Acids Res* 2006, **34**:W435-439.

- 2
3 37. Cai Y, Gonzalez JV, Liu Z, Huang T: **Computational systems biology methods in**
4 **molecular biology, chemistry biology, molecular biomedicine, and biopharmacy.**
5 *Biomed Res Int* 2014, **2014**:746814.
6
7
8
9 38. Majoros WH, Pertea M, Salzberg SL: **TigrScan and GlimmerHMM: two open source**
10 **ab initio eukaryotic gene-finders.** *Bioinformatics* 2004, **20**:2878-2879.
11
12 39. Elsik CG, Mackey AJ, Reese JT, Milshina NV, Roos DS, Weinstock GM: **Creating a**
13 **honey bee consensus gene set.** *Genome Biol* 2007, **8**:R13.
14
15
16 40. Lowe TM, Eddy SR: **tRNAscan-SE: a program for improved detection of transfer**
17 **RNA genes in genomic sequence.** *Nucleic Acids Res* 1997, **25**:955-964.
18
19
20 41. Gardner PP, Daub J, Tate J, Moore BL, Osuch IH, Griffiths-Jones S, Finn RD, Nawrocki
21 EP, Kolbe DL, Eddy SR, Bateman A: **Rfam: Wikipedia, clans and the "decimal"**
22 **release.** *Nucleic Acids Res* 2011, **39**:D141-145.
23
24
25 42. UniProt C: **UniProt: a hub for protein information.** *Nucleic Acids Res* 2015,
26 **43**:D204-212.
27
28
29 43. Kanehisa M, Goto S, Sato Y, Kawashima M, Furumichi M, Tanabe M: **Data, information,**
30 **knowledge and principle: back to metabolism in KEGG.** *Nucleic Acids Res* 2014,
31 **42**:D199-205.
32
33
34 44. Edgar RC: **MUSCLE: a multiple sequence alignment method with reduced time and**
35 **space complexity.** *BMC Bioinformatics* 2004, **5**:113.
36
37
38 45. Talavera G, Castresana J: **Improvement of phylogenies after removing divergent and**
39 **ambiguously aligned blocks from protein sequence alignments.** *Syst Biol* 2007,
40 **56**:564-577.
41
42
43 46. Posada D: **Using MODELTEST and PAUP* to select a model of nucleotide**
44 **substitution.** *Curr Protoc Bioinformatics* 2003, **Chapter 6**:Unit 6 5.
45
46
47 47. Huelsenbeck JP, Ronquist F: **MRBAYES: Bayesian inference of phylogenetic trees.**
48 *Bioinformatics* 2001, **17**:754-755.
49
50
51 48. Stamatakis A: **RAxML-VI-HPC: maximum likelihood-based phylogenetic analyses**
52 **with thousands of taxa and mixed models.** *Bioinformatics* 2006, **22**:2688-2690.
53
54
55 49. Tamura K, Peterson D, Peterson N, Stecher G, Nei M, Kumar S: **MEGA5: molecular**
56
57
58
59
60
61
62
63
64
65

1 evolutionary genetics analysis using maximum likelihood, evolutionary distance, and
2 maximum parsimony methods. *Mol Biol Evol* 2011, **28**:2731-2739.

- 3
4 50. Wang MS, Yang HC, Otecko NO, Wu DD, Zhang YP: **Olfactory genes in Tibetan wild**
5 **boar**. *Nat Genet* 2016, **48**:972-973.
6
7 51. Yang Z: **PAML 4: phylogenetic analysis by maximum likelihood**. *Mol Biol Evol* 2007,
8 **24**:1586-1591.
9
10 52. Chen H: **Population genetic studies in the genomic sequencing era**. *Zool Res* 2015, **36**:
11 223-232.
12
13 53. Wu, D, D; Wang, M, S; Zeng, Y; Wang, X; Nie, W, H; Wang, J, H; Su, W, T; Otecko, N,
14 O; Xiong, Z, J; Wang, S; Qu, K, X; Wang, W; Dong, Y; Zhang, Y, P (2017): Supporting
15 data for "Draft genome of the Gayal, *Bos frontalis*" GigaScience Database.
16 <http://dx.doi.org/10.5524/100354>
17
18
19
20
21
22
23
24
25
26
27
28
29
30

31 **Figure Legends:**

- 32
33 Figure 1. A picture showing a female gayal (*Bos frontalis*, provided by Kai-Xing Qu).
34
35 Figure 2. Karyotype of the gayal used for genome sequencing (provided by Wen-Hui
36 Nie).
37
38 Figure 3. 17-mer frequency distribution of sequencing reads.
39
40 Figure 4. Phylogenetic trees of gayal and other bovine species. (A) Tree constructed
41 based on maximum likelihood method, (B) Tree constructed using Bayesian
42 inference.
43
44 Figure 5. Maximum likelihood trees of gayal and other bovine species using whole
45 complete mtDNA. IDs in parentheses are GenBank accession number.
46
47 Figure 6. Divergence time estimated between gayal and other bovine species.
48
49
50
51
52
53
54
55

56 **Tables:**

- 57
58 Table 1. Statistics of the completeness of the hybrid *de novo* assembly of *Bos frontalis*
59 genome
60
61
62
63
64
65

Terms	Contig		Scaffold	
	Size	number	Size	number
N90	2,461	211577	158,610	1357
N80	5,335	140237	1,060,177	800
N70	8,109	99930	1,668,147	587
N60	11,044	71764	2,170,469	437
N50	14,405	50585	2,737,757	320
Max length	208,099		13,764,521	
Total length	2,669,378,334		2,848,570,279	
Total number		583373		460,059
Average length	4575		6,191	
Number>=500bp		394757		116481
Number>=1000bp		300178		53989
Number>=2000bp		229796		19915
Number>=5000bp		146493		5387

Table 2. Statistics of the completeness of the assembled genomes for *Bos frontalis* and close related species by BUSCO (version 2)

Species	Terms	Complete(C)	Complete and single-copy (S)	Complete and duplicated (D)	Fragmented (F)	Missing (M)
gayal	Number	3494	3434	60	319	291
	Proportion	85.14%	83.67%	1.46%	7.77%	7.09%
zebu	Number	3698	3644	54	158	248
	Proportion	90.11%	88.79%	1.32%	3.85%	6.04%
wisent	Number	3794	3763	31	180	130
	Proportion	92.45%	91.69%	0.76%	4.39%	3.17%
yak	Number	3841	3809	32	138	125
	Proportion	93.59%	92.81%	0.78%	3.36%	3.05%
buffalo	Number	3817	3780	37	142	145
	Proportion	93.01%	92.11%	0.90%	3.46%	3.53%
bison	Number	3779	3735	44	165	160
	Proportion	92.08%	91.01%	1.07%	4.02%	3.90%

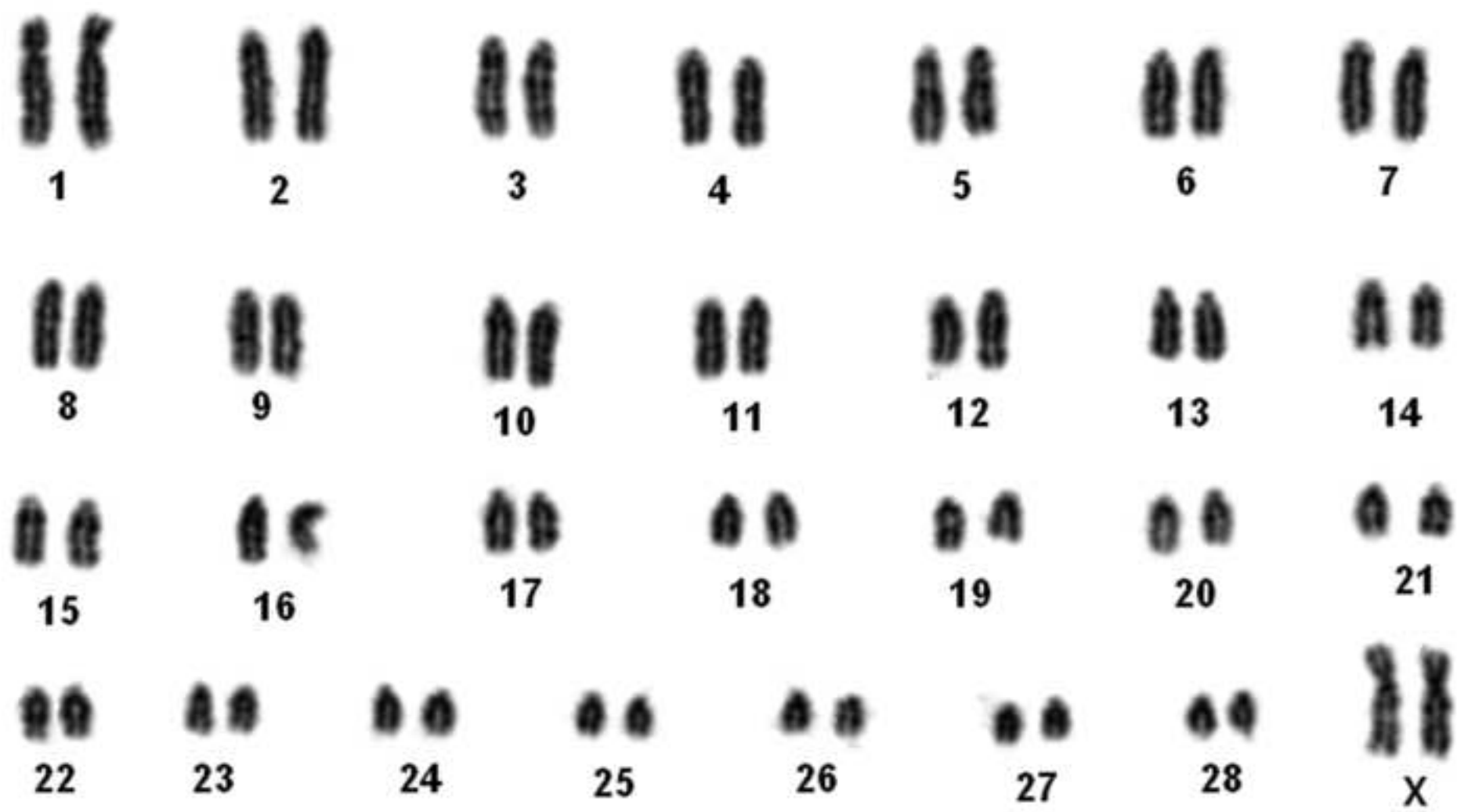
Table 3. Statistics of repeats in *Bos frontalis* genome.

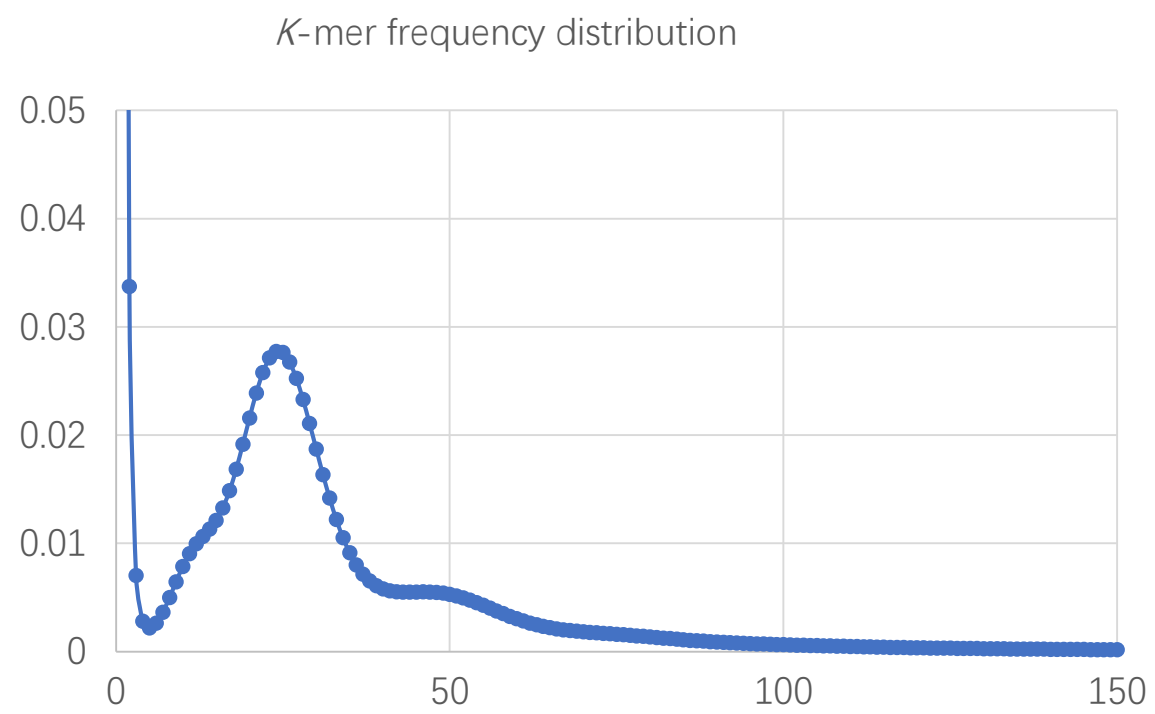
Type	Repeat Size (bp)	% of genome
Trf	17,696,175	0.62
Repeatmasker	868,885,926	30.50
Proteinmask	265,003,148	9.30
<i>De novo</i>	917,371,710	32.20
Total	1,371,023,312	48.13

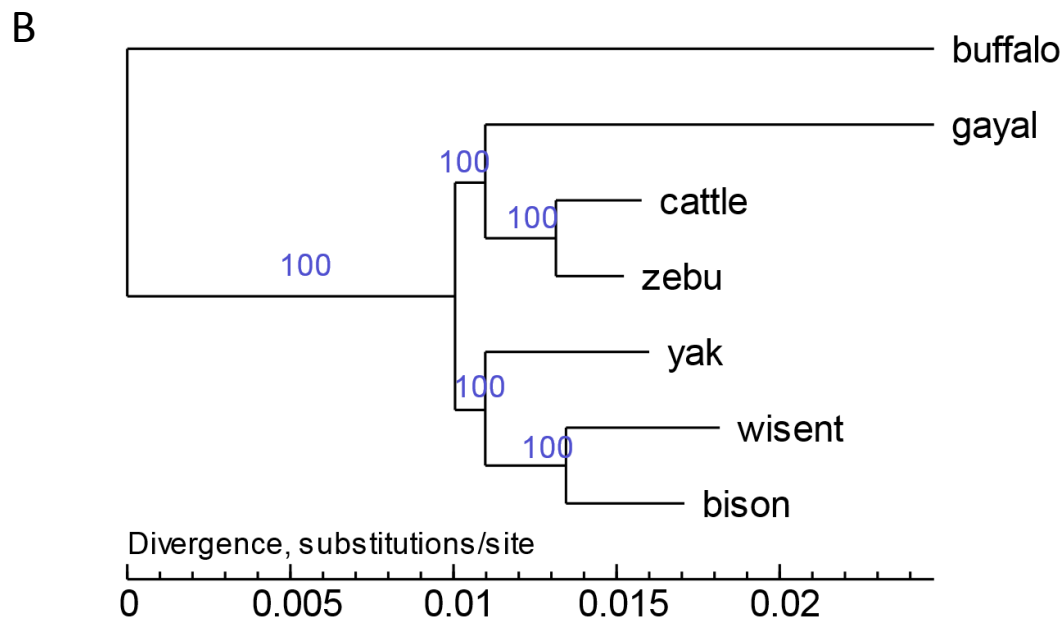
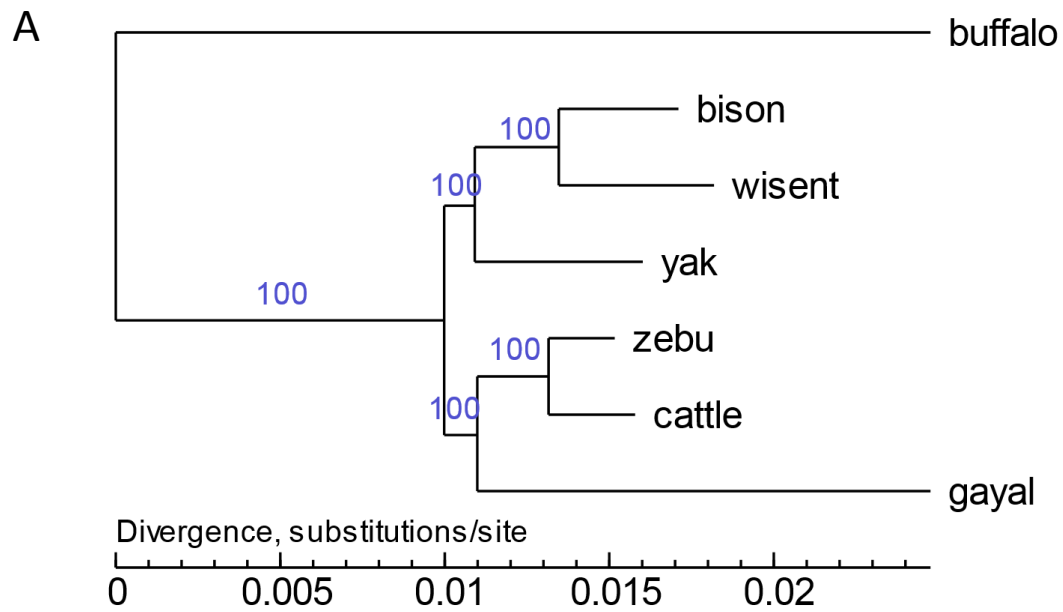
Table 4. General statistics of predicted protein-coding genes.

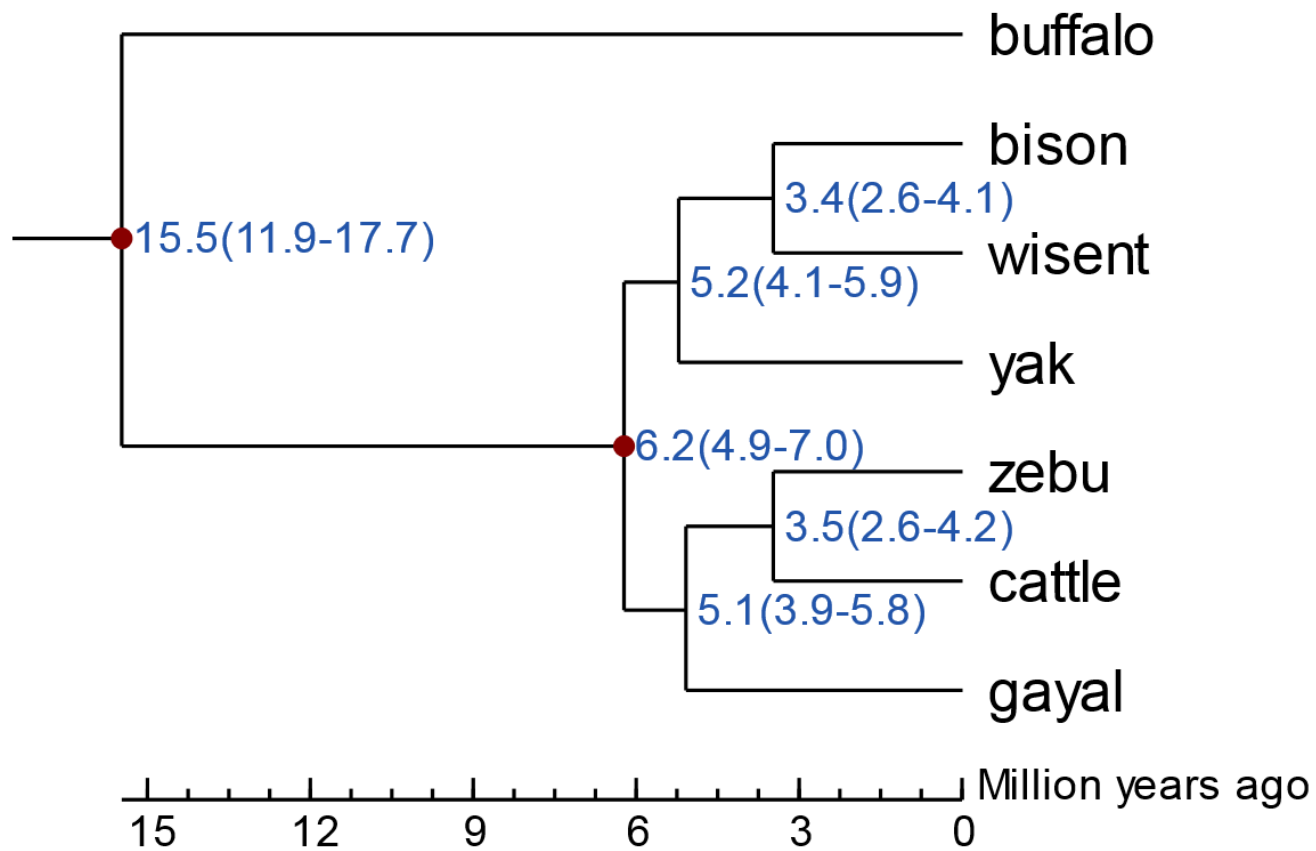
Gene set		Total	Exon number	CDS length (bp)	mRNA length (bp)	Exons per gene	Exon length (bp)	Intron length (bp)
Homolog	<i>Bos taurus</i>	19,666	141,323	1,325	20,618	7.19	184	3,118
	<i>Canis familiaris</i>	17,627	121,986	1,323	20,802	6.92	191	3,290
	<i>Homo sapiens</i>	24,783	146,172	1,108	17,567	5.89	187	3,360
	<i>Sus scrofa</i>	20,283	121,282	1,142	16,288	5.97	191	3,041
	<i>Rattus norvegicus</i>	17,988	117,965	1,277	19,469	6.55	194	3,273
	<i>Ovis aries</i>	20,947	147,367	1,287	20,973	7.03	183	3,261
<i>De novo</i>	AUGUSTUS	41,227	180,664	1,127	22,786	4.38	257	6,403
	GlimmerHMM	27,067	104,294	874	5,433	3.85	226	1,597
	Genescan	46,598	297,828	1,321	36,828	6.39	206	6,585
Glean (final)		26,667	87,392	1,156	4,996	3.27	352	1,686














Click here to access/download
Supplementary Material
SupplementaryInfor.doc



23 August 2017

To
The Editorial Office,
GigaScience.

Dear Editor,

RE: Draft genome of the gayal, *Bos frontalis*.

We are very grateful for the constructive review and the valuable suggestions received on our paper titled above (GIGA-D-17-00116). We have carefully considered all the comments and suggestions. It is with much pleasure that we re-submit the revised manuscript for your consideration for publication in your highly reputable journal.

We have carefully revised our citations and references, sample origin, and the novelty of our research. We have also re-estimated the genome size using K-mer ratio using only the reads that have passed quality filtering. For a detailed description of all the changes in the manuscript, please find a separate point-by-point response to each of the comments raised by the editor and the reviewers.

We believe the review has greatly improved our manuscript and that our revisions sufficiently address all the review comments. Thank you in advance for considering this manuscript for publication.

Yours sincerely,
Dong-Dong Wu and Co-authors