1    **The First High-Quality Reference Genome of Sika Deer Provides**

2    **Insights for High-Tannin Adaptation**

3

4    Xiumei Xing[*,#,1], Cheng Ai[#,2], Tianjiao Wang[#,1], Yang Li[#,1], Huitao Liu[#,1], Pengfei Hu[#,1],

5    Guiwu Wang[#,1], Huamiao Liu[1], Hongliang Wang[1], Ranran Zhang[1], Junjun Zheng[1],

6    Xiaobo Wang[2], Lei Wang[1], Yuxiao Chang[2], Qian Qian[2], Jinghua Yu[3], Lixin Tang[1],

7    Shigang Wu[2], Xiujuan Shao[2], Alun Li[2], Peng Cui[2], Wei Zhan[4], Sheng Zhao[2], Zhichao

8    Wu[2], Xiqun Shao[1], Yimeng Dong[1], Min Rong[1], Yihong Tan[3], Xuezhe Cui[1], Shuzhuo

9    Chang[1], Xingchao Song[1], Tongao Yang[1], Limin Sun[1], Yan Ju[1], Pei Zhao[1], Huanhuan

10   Fan[1], Ying Liu[1], Xinhui Wang[1], Wanyun Yang[1], Min Yang[1], Tao Wei[1], Shanshan

11   Song[1], Jiaping Xu[1], Zhigang Yue[1], Qiqi Liang[*,5], Chunyi Li[*,1], Jue Ruan[*,2], Fuhe

12   Yang[*,1]

13

14   *[1] Key Laboratory of Genetics, Breeding and Reproduction of Special Economic Animals,*

15   *Ministry of Agriculture, Institute of Special Animal and Plant Sciences, Chinese*

16   *Academy of Agricultural Sciences, Changchun 130112, China*

17   *[2] Guangdong Laboratory for Lingnan Modern Agriculture, Genome Analysis*

18   *Laboratory of the Ministry of Agriculture, Agricultural Genomics Institute at Shenzhen,*

19   *Chinese Academy of Agricultural Sciences, Shenzhen 518120, China*

20   *[3] CAS Key Laboratory of Forest Ecology and Management, Institute of Applied Ecology,*

21   *Chinese Academy of Sciences, Shenyang 110016, China*

22   *[4] Annoroad Gene Technology (Beijing) Co., Ltd, Beijing Economic-Technological*

23   *Development Area, Beijing 100176, China*

24   *[5] Novogene Bioinformatics Institute, Chaoyang District, Beijing 100083, China*

25

26   [#] These authors contributed equally.

27

28   * Corresponding authors

29     Xiumei Xing, Qiqi Liang, Chunyi Li, Jue Ruan, Fuhe Yang

30     E-mail: xingxiumei@caas.cn (Xing X), liangqiqi@novogene.com (Liang Q),

31     lichunyi1959@163.com (Li C), ruanjue@caas.cn (Ruan J), yangfh@126.com (Yang F).

32

33     **Running title**: Xing *X et al / Chromosome-Level Reference Genome Sequence of Sika*

34     *Deer*

35

36     **KEYWORDS:** Sika deer; Whole-genome sequencing; Chromosome-scale assembly;

37     Oak leaves; Tannin tolerance

38

39     **Total number of words**

40     7489

41     **Total number of figures**

42     **3**

43     **Total number of tables**

44     **1**

45     **Total number of supplementary figures**

46     **17**

47     **Total number of supplementary tables**

48     **19**

49

## Abstract

Sika deer are known to prefer oak leaves, which are rich in tannins and toxic to most mammals; however, the genetic mechanisms underlying their unique ability to adapt to living in the jungle are still unclear. In identifying the mechanism responsible for the tolerance of a highly toxic diet, we have made a major advancement in the elucidation of the genomics of sika deer. We generated the first high-quality, chromosome-level genome assembly of sika deer and measured the correlation between tannin intake and RNA expression in 15 tissues through 180 experiments. Comparative genome analyses showed that the *UGT* and *CYP* gene families are functionally involved in the adaptation of sika deer to high-tannin food, especially the expansion of *UGT* genes in a subfamily. The first chromosome-level assembly and genetic characterization of the tolerance toa highly toxic diet suggest that the sika deer genome will serve as an essential resource for understanding evolutionary events and tannin adaptation. Our study provides a paradigm of comparative expressive genomics that can be applied to the study of unique biological features in non-model animals.

## Introduction

66

67    Cervidae consists of 55 extant deer species and constitutes the second largest family in

68    terrestrial artiodactyls. Sika deer (*Cervus nippon*) is naturally distributed throughout

69    East Asia and is one of the best-known deer species producing velvet antlers [1,2], a

70    valuable ingredient in traditional Chinese medicine [3]. Among other deer species [4-

71    6], sika deer has unique characteristics, such as a geographic distribution that is

72    significantly more coincident with oak trees (Figure 1A) and an ability to tolerate a

73    high-tannin diet, mainly consisting of oak leaves. Notably, oak leaves, which are rich

74    in tannins and toxic to most mammals, such as cattle, which are related to sika deer [7],

75    are conversely found to increase the reproductive rate and fawn survival rate of sika

76    deer. Thus, oak leaves are essential for maintaining healthy sika deer in wild and farmed

77    populations. Some studies have concluded that tannins are not toxic to sika deer because

78    of the rumen microbes and fermentation patterns of these deer [8]. However, knowledge

79    is scarce regarding the genetics and mechanism underlying the ability to detoxify a

80    high-tannin diet.

81    Whole-genome sequencing has become a more popular technology with which to

82    explore the taxonomy, evolution and biological phenomena of organisms at the

83    molecular level [9], compared with morphological, histological and other analyses [10-

84    12]. For example, a series of studies investigated the genomes of 11 deer and 33 other

85    ruminant species and identified some genes that are involved in ruminant headgear

86    formation, rapid antler regeneration, and reindeer adaptation to the long days and nights

87    in the Arctic region [6,13,14]. The chromosome-level reference genome for sika deer

88    is in high demand compared with that for other ruminants such as bovines [15,16], and

89    it will provide novel genomic and molecular evolutionary information on the

90    exceptional characteristics of the sika deer.

91    Here, we report the chromosome-level genome assembly of a female sika deer, as

92    well as the RNA sequencing of 15 tissue types in sika deer treated with 3 levels of a

93    high-tannin diet. The findings provide important resources to help elucidate the genetic

94   mechanisms underlying the high-tannin food tolerance of sika deer. Our high-quality

95   sika deer genome will be of great importance to researchers who study the common

96   characteristics of deer and other ruminants and could even serve as a reference deer

97   genome. The well-designed RNA expression experiments used in this study also

98   provide a paradigm for studying novel features in nonmodel animals.

99

100

101   **Results**

102   **De novo assembly of a Cervus nippon reference genome**

103   We collected DNA from a female sika deer (*Cervus nippon*) and identified a total of 66

104   chromosomes, including 64 autosomes and one pair of sex chromosomes (XX)

105   (Additional file 1: Figure S1). A large set of data was acquired for assembly using a

106   combination of four technologies. 1) A total of 242.9 Gb of clean data (~93.4×) were

107   obtained from paired-end sequencing (Illumina HiSeq), with the genome size (2.6 Gb)

108   estimated by the 25 K-mer distribution (Additional file 2: Table S1 and Additional file

109   1: Figure S2). 2) A total of 150.4 Gb (~57.7×) of PacBio RSII long reads (single-

110   molecule real-time sequencing) were also acquired (Additional file 2: Table S2). The

111   wtdbg2 [17] assembler yielded 2,040 primary contigs using PacBio reads with a contig

112   N50 size of 23.6 Mb and the longest at 93.6 Mb (Additional file 2: Table S3). These

113   contigs were then polished using the Quiver algorithm [18] with default parameters.

114   Genome-wide base-level correction was performed using Illumina short reads aligned

115   to the published genome with BWA (v0.7.10-r943-dirty), and inconsistencies between

116   the genome and the reads were identified with SAMTools/VCFtools (v1.3.1). These

117   inconsistencies were corrected by our in-house script to produce a highly accurate

118   assembly. 3) The previous contigs were clustered into chromosome-scale scaffolds

119   using high-throughput chromosome conformation capture (Hi-C) proximity-guided

120   assembly (Figure 1B) to produce the final reference assembly, named MHL_v1.0,

121   totaling 2.5 Gb of sequence with a contig N50 of 23.6 Mb and a scaffold N50 of 78.8

122    Mb (**Table 1**). The resulting assembly contained 2,481,763,803 bp reliably anchored

123    on chromosomes, accounting for 99.24% of the whole genome (Additional file 2: Table

124    S4). 4) A total of 264 Gb of optical mapping (using BioNano Genomics Irys) data were

125    also used to generate *de novo*-assembled optical maps with a scaffold N50 of 1.974 Mb,

126    which was sequentially compared with MHL_v1.0 to identify the misoriented contigs

127    and improve the final validated reference assembly (Additional file 1: Figure S3).

128        To validate our assembly, MHL_v1.0 was compared with the previously published

129    red deer [19] genome (Additional file 1: Figure S4). Both the inconsistency of the

130    synteny analysis and the improper density of Hi-C proximity maps identified 34

131    inaccurate junctions, which were considered potential inversions and misassemblies

132    (Additional file 1: Figures S4 and S5). The aforementioned optical maps were used to

133    determine whether the 34 inaccurate junctions were breakpoints or new joint regions

134    after the replacement. We found that 10 inaccurate junctions were supported by the

135    optical maps, and those junctions were then manually inspected and correlated.

136    Additionally, another 142 potential misjoined contigs were found by comparing our

137    MHL_v1.0 assembly with the optical maps. The paired-end Illumina short reads were

138    then mapped to the final assembly, and all 142 disagreements were checked manually

139    and found to be sequential in the comparison results. We further compared MHL_v1.0

140    with the twenty published genomes of Cervidae, including red deer (*Cervus elaphus*)

141    and reindeer (*Rangifer tarandus*). The results showed that the scaffold N50 length and

142    ungapped sequence length of the MHL_v1.0 assembly were greater than those

143    previously published (Additional file 2: Table S5). We compared three other

144    chromosome-level ruminant genomes (cattle, goat, and red deer) with MHL_v1.0.

145    Multiple chromosome fission/separation events were detected among these four

146    genomes, and we found that the sika deer genome had the highest chromosome

147    collinearity with red deer (Figure 1C and Additional file 1: Figure S6).

148        Finally, we downloaded a total of 2,715 EST sequences belonging to sika deer from

149    the NCBI dbEST database and aligned them against MHL_v1.0. We found that 95.95%

150     of the EST sequences (coverage rate > 90%) matched our sika deer genome MHL_v1.0.

151     Evaluation of our MHL_v1.0 using CEGMA software     showed that 97.18% of the full

152     length of 248 genes in the core gene set was predicted. Benchmarking Universal Single-

153     Copy Orthologs (BUSCO) analysis of the gene set showed that complete BUSCO

154     accounted for 3,880 (of 4,104; 94.60%) genes, which is better than the results obtained

155     for the water buffalo (*Bubalus bubalis*, 93.6%) [12] and domestic goat (*Capra hircus*,

156     *82%*) [20]. After aligning Illumina short reads (93.4×) against MHL_v1.0, the base

157     level error rate was estimated to be 1.1e-5 (Additional file 2: Table S6).

158     **Genome annotation**

159     Homology and *de novo* repetitive sequence annotation results showed that repetitive

160     sequences accounted for approximately 45.38% of MHL_v1.0, which is consistent with

161     the percentages published for other mammals (Additional file 2: Tables S7 and S8),

162     including humans (44.8%) [21], water buffalo (45.33%) [12] and sheep (42.67%) [22].

163     As in other published mammalian genomes, long interspersed nuclear elements

164     (LINEs), short interspersed nuclear elements (SINEs) and long terminal repeats (LTRs)

165     were also the most abundant elements in MHL_v1.0 (29.56%, 7.63% and 5.38% of the

166     total number of elements, respectively) (Additional file 1: Figure S7). The main features

167     of MHL_v1.0 are summarized and shown in Additional file 1: Figure S8.

168       A total of 21,449 protein-coding genes were predicted using the combined methods

169     of homology and *de novo* annotations with transcriptome data (mapping rate of 93.43%

170     for 1.2 billion RNA-Seq reads), and 90.1% of the protein-coding genes were

171     functionally annotated (Additional file 2: Table S9). The average coding sequence

172     (CDS) length per gene was 1,617 bp, the exon number per gene was 9.29, and the

173     average length per exon was 174 bp; these values are similar to those in other mammals

174     (Additional file 2: Table S10). To verify the accuracy of our gene predictions and to

175     assess the annotation completeness of MHL_v1.0, we checked core gene statistics using

176     the BUSCO software. A total of 3,907 (of 4,104; 95.20%) (Additional file 2: Table S11)

177     highly conserved core proteins in mammals were recovered from our predictions.

**Analyses of phylogeny and demographics**

A phylogenetic tree (**Figure 2A**) based on 19 mammals spanning the orders Primates, Rodentia, Artiodactyla and Cetacea was constructed with the maximum-likelihood method using 748 identified single-copy orthologous genes. The results showed that sika deer was in the same clade as red deer (Figure 2A), which is consistent with the cladistic data [23]. The divergence time between sika deer and red deer was estimated to be approximately 2.5 million years ago (MYA) (Figure 2A and Additional file 1: Figure S9).

To examine the changes in effective population size (Ne) of the ancestral populations, a Pairwise Sequential Markovian Coalescent (PSMC) analysis was applied to sika deer, cattle [16] and buffalo [12] (Figure 2B). Demographic analysis showed that the Ne of the sika deer sharply declined during the two large glaciations: the Qingzang movement (1.7-3.6 MYA) and Penultimate Glaciation (0.13-0.3 MYA), and the sika deer underwent a long period of population bottlenecks. Subsequently, the Ne increased greatly after that period, suggesting that these deer had adapted to the specific habitat, probably due to the monsoon climate in East Asia. During the same period, the populations of cattle and buffalo recovered soon after a decline and shrank again. During Marine Isotope Stage 4 (0.058-0.074 MYA) and the last glacial maximum (LGM, ~0.02 MYA), sika deer suffered population bottlenecks again (Figure 2B), which may also be the reason modern sika deer populations have very low genetic diversity [23].

**Gene family evolution**

We identified a total of 9,830 homologous gene families in MHL_v1.0 by comparing the predicted protein sequences of sika deer with those of 19 mammals spanning the orders Primates, Rodentia, Artiodactyla and Cetacea (Additional file 2: Table S12 and Additional file 1: Figure S10).

Based on the hypothesis that potential genomic adaptations are related to genes that are under positive selection in the sika deer lineages [24], we identified 55 positively

206  selected genes (PSGs), which were calculated using the branch-site models and

207  validated using likelihood ratio tests (Additional file 2: Table S13). The PSGs were

208  found to be involved in the PI3K-Akt signaling pathway (ko04151), VEGF signaling

209  pathway (ko04370) and pathways in cancer (ko05200), among others. These pathways

210  were reportedly related to antler growth [25,26].

211      The number of genes in a gene family has been proposed as a major factor

212  underlying the adaptive divergence of closely related species. To depict the gene family

213  evolution, we identified 972 significantly contracted and 879 significantly expanded

214  gene families in sika deer compared with other species (Figure 2A). The expanded gene

215  families were mainly enriched in the signal transduction pathways of environmental

216  perception (olfactory transduction, G protein-coupled receptors, neuroactive ligand-

217  receptor interaction, corrected $P$-value < 0.05), enzymatic activity (transferase activity,

218  transferring hexosyl groups, carboxypeptidase activity and L-lactate dehydrogenase

219  activity, corrected $P$-value < 0.05), feeding behavior (salivary secretion,

220  neurotransmitter secretion, corrected $P$-value < 0.05) and drug metabolism (drug

221  metabolism - other enzymes, drug metabolism - cytochrome P450, metabolism of

222  xenobiotics by cytochrome P450, corrected $P$-value < 0.05) (Additional file 2: Tables

223  S14 and S15). The contracted gene families were mainly related to lipid metabolism

224  pathways (linoleic acid metabolism and ether lipid metabolism, corrected $P$-value <

225  0.05), ion transportation (calcium ion binding, anion transport, and iron ion binding,

226  corrected $P$-value < 0.05) and regulation of basic biological processes (regulation of

227  developmental and apoptotic processes, corrected $P$-value < 0.05) (Additional file 2:

228  Tables S16 and S17).

229  **Exceptional expansion of the UGT gene family in the sika deer genome**

230  Gene gains and losses are one of the primary contributors to functional changes. To

231  better understand the evolutionary dynamics of genes, we assessed the expansion and

232  contraction of the gene ortholog clusters among 19 species. The uridine 5'-diphospho-

233  glucuronosyltransferase (UDP-glucuronosyltransferase, *UGT*) gene families were at

234 the top 27 of 879 significantly expanded gene families, which have been reported to

235 play a role in the catabolism of exogenous compounds [27-29]. Phylogenetic analysis

236 revealed that the 257 *UGT* genes could be classified into 7 lineages (Figure 3A and

237 Additional file 1: Figure S11), while in the sika deer genome, we found two lineage-

238 specific monophyletic expansions of the *UGT2B* and *UGT2C* subfamilies (Figure 3B).

239 In the *UGT2B* subfamily, 15 copies were found in the sika deer genome, which was

240 more than that in any other species assessed in this study (Additional file 2: Table S18).

241 Sika deer had relatively lower expanded gene numbers in the *UGT2C* subfamily than

242 in the *UGT2B* subfamily (Additional file 2: Table S18). Taken together, these results

243 prompt us to propose that the exceptional expansion of the *UGT* gene family may be

244 the key genetic basis for the tolerance of high-tannin food, namely, oak leaves, by the

245 sika deer.

246 **Transcriptomic analysis of 15 tissues of sika deer treated with a high-tannin diet**

247 Sika deer adapted well to living in the forest and have consumed a high-tannin diet of

248 Mongolian oak (*Quercus mongolica*) leaves (MOL) for a long time; whether the

249 underlying genetic adaptation and molecular mechanism are associated with the special

250 expansion of *UGT* gene families is an interesting question. We used 9 deer fawns to

251 conduct a feeding trial with different tannin-containing (0%, 50%, 100%) diets, and 3

252 mature deer (100%) were used as a comparison group. Transcriptome sequencing was

253 performed on 15 tissues of all experimental individuals (Additional file 2: Table S19).

254 A total of 1.44 Tb of transcriptional data from 180 samples were obtained using the

255 Illumina platform, and the 17,233 differentially expressed genes (DEGs) were analyzed

256 by pairwise comparison of each group (Additional file 1: Figure S12). The liver is the

257 major organ associated with *UGT* activity, and *UGT* expression was highest in the liver

258 among the fifteen tissues examined (Figure 3C). Although *UGT* genes were also highly

259 expressed in the liver tissue of cattle, they did not respond to high MOL levels

260 (Additional file 1: Figure S13). We compared different MOL levels in sika deer and

261 identified 3,222 and 15 DEGs in liver and kidney tissue, respectively.

262　　　　After inspecting all the expanded/contracted gene families and DEGs in liver tissue,

263　29 genes were found to play roles in the P450 pathway. Of these, 20 were expanded

264　genes, 12 were DEGs, and 3 were contracted genes. The interaction network of these

265　genes is shown in Figure 3D. Among these key genes, *UGT2B4* and *UGT2B31* were

266　both significantly upregulated in high-tannin liver tissue and expanded in the sika deer

267　genome. Therefore, we hypothesized that *UGT2B4* and *UGT2B31* are major genes in

268　sika deer with high-tannin adaptation.

269　　　　Interestingly, in liver tissue, tannins can drive the expression of many *UGT* genes

270　in a dose-dependent manner. Overall, when compared among different MOL levels and

271　ages (y0, y50, y100 and m100), eight differentially expressed *UGT* genes were

272　discovered, among which two were downregulated genes of the *UGT3A* subfamily and

273　six were upregulated genes in the *UGT2B* and *UGT2C* subfamilies (Figure 3E).

274　Furthermore, we found that all of these upregulated *UGT* genes in the liver were located

275　on sika deer chromosome 27 (Figure 3F). With the increase in tannin content intake,

276　the *UGT3A* subfamily genes in the liver were inhibited; nevertheless, *UGT* gene copies

277　in the *UGT2B* and *UGT2C* families were increased, suggesting that the response of

278　*UGT* gene expression to tannin was mainly upregulated. Moreover, in the kidney tissue,

279　two DEGs belonged to the *UGT2C* family. Five differentially expressed *CYP* genes

280　were upregulated, whereas gene families encoding *GST* and *SULT* were all

281　downregulated after the deer were fed a high-tannin diet. According to previous studies,

282　sika deer share common pathways with koala, including the drug metabolism-

283　cytochrome P450 signal pathway [11]. The detoxification genes in sika deer showed

284　opposite expression patterns compared with the genes in koala [11] (Additional file 1:

285　Figure S14). These results indicate that sika deer may utilize a different adaptive

286　strategy from that of koala to survive on a diet of highly toxic food.

287　**Ability to tolerate a high-tannin diet**

288　The sika deer diet of MOL contains high levels of tannins that would be lethal to most

289　other mammals. The main detoxification reactions are traditionally categorized into

290    phase I and phase II reactions. Currently available evidence indicates that among these,

291    the *CYP*, *UGT*, *GST*, and *SULT* gene families have the greatest importance in

292    xenobiotic metabolism. Based on the aforementioned mechanism, genes involved in

293    those pathways were examined using gene family and transcriptome analyses.

294    A total of 13 DEGs were detected from the *CYP2* family in sika deer liver, but only

295    5 were differentially expressed with increasing tannin contents in the diet. Five *GST*

296    genes and 3 *SULT* genes were found to be differentially expressed in the liver, but all

297    were downregulated with increasing tannin contents in the diet.

298    The functional importance of these *UGT* genes was further investigated through

299    analysis of their expression levels in sika deer, showing that they had particularly high

300    expression in the liver tissue, which is consistent with their role in detoxification. The

301    mechanism of the glucuronidation reaction is that *UGT* enzymes catalyze the transfer

302    of the glucuronosyl group from uridine 5'-diphospho-glucuronic acid (UDPGA) to the

303    tannin molecules, generating the glucuronidated metabolite, which is more polar and

304    more easily excreted than the tannin molecule (Figure 3F). Most of these expressed

305    *UGT* genes belonged to *UGT2B*. These phenotypes suggest that *UGT* genes in *UGT2B*

306    have an important role in detoxification; the upregulated expansion of *UGT* genes

307    would result in higher enzyme levels, which would enhance the ability of sika deer to

308    detoxify the high-tannin diet.

309    Among the genes related to the metabolism of drugs and exogenous substances,

310    *UGT* and *CYP* genes were found to be functionally involved in detoxification,

311    especially *UGT* genes in the *UGT2B* family. In short, these findings imply that the

312    unique expansion of the *UGT* gene family is mainly responsible for the toleration of

313    high-tannin food, namely, oak leaves, by sika deer (Additional file 1: Figure S15).

314

## Discussion

316    Cervidae is the second largest family in Artiodactyla [30] and has significant scientific

317    [1] and economic [3] value. Although several other deer genomes have recently been

318   reported, the lack of high-quality genome sequences of sika deer, one of the novel

319   species in the family, has hindered the elucidation of the molecular mechanisms

320   underlying important distinct biological characteristics of sika deer, such as the full

321   regeneration of the antlers. Here, we sequenced the genome of sika deer and assembled

322   it at the chromosome level using combined technologies of SMRT, Illumina sequencing

323   and Hi-C. The high percentage and accuracy rate of the genome structure, base calling,

324   gene set validation and quality of gene annotation demonstrated that our assembled sika

325   deer genome was of high quality and could be effectively used as a reference genome

326   for deer species.

327      The geographic distribution of sika deer is highly coincident with that of oak, and

328   sika deer have a preference for grazing on high-tannin oak leaves [31], suggesting that

329   this adaptation may be a positive selection during evolution. In terms of food adaption,

330   sika deer are not unique. For example, pandas, dogs and koalas have also undergone

331   adaptive food evolution; pandas can eat bamboo despite being carnivorous [32], dogs

332   can adapt to a diet of starchy foods [33], and koalas can eat toxic eucalyptus leaves [11].

333   Divergent adaptive pathways and related genes are known to be involved in this

334   adaptation. In this study, we found that among the genes related to toxin degradation,

335   only those from the *UGT* gene family [34], especially the *UGT2B* family, were

336   significantly expanded. Furthermore, transcriptomic studies showed that *UGT* gene

337   expression was strongly correlated with the quantity of tannin intake, i.e., it was dose

338   dependent. The expression of specific extended gene copies in the *UGT2B* family was

339   prominently increased after the tannin feeding treatment. These results suggest that

340   genes in the *UGT* family, especially in the *UGT2B* subfamily, are associated with the

341   adaptation of sika deer to a high-tannin diet.

342      It is generally believed that rumen microorganisms play a role in the digestion of

343   tannins [35,36]. However, as other ruminants, such as cattle and sheep, are not well

344   adapted to high-tannin diets (Additional file 1: Figure S16), we speculate that during a

345   long period of coexistence with oak trees during evolution, sika deer may have

346     developed genetic adaptive mechanisms. As expected, we found evidence for this

347     phenomenon at the genome level through high-quality sequencing. Transcriptomic

348     results also revealed that changes in gene expression were involved in Na and K ion

349     channels. The Na and K balance (water and salt metabolism) is essential for the basic

350     metabolism of organisms. These genetic responses have enabled sika deer to adapt to

351     oak leaves as an advantageous rather than a hazardous material for consumption.

352

353

354     **Conclusion**

355     The sika deer genome assembled in this study provides, to our knowledge, the highest

356     quality deer genome to date. The comprehensive characterization of the sika deer

357     genome along with the transcriptomic data presented herein provides a framework used

358     to elucidate its evolutionary events, revealing the mechanism of the unique attributes

359     and tannin adaptation. Through detailed genomics and transcriptomics analyses, we

360     identified the most likely mechanism of tannin degradation in sika deer. We also

361     depicted possible molecular mechanisms for the jungle adaptability of deer, and the

362     methodologies we used in this study will also provide a reference for the study of the

363     adaptation mechanism of animals to "toxic" foods. Chromosome-scale assembly of sika

364     deer genomes could be used for many applications, including the study of structural

365     variations in large genomic regions, expected recombination frequencies in specific

366     genomic regions, target sequence characterization and modification for gene editing.

367     Moreover, this study provides a valuable genomic resource for research on the genetic

368     basis of sika deer's distinctive physiological features, such as the full regeneration of

369     deer antlers, and on Cervidae genome evolution. Our study also contributes to

370     conservation and utilization efforts for this antler-growing species.

371

372

373     **Materials and methods**

374 **Method details**

375 *Sampling preparation*

376 A female sika deer (*Cervus nippon*) from Jilin Province was used for *de novo* genome

377 sequencing. DNA was extracted from whole blood with a BioTeke DP1102 kit (solution)

378 according to the manufacturer's instructions. After slaughtering the experimental

379 animals, tissue sampling was carried out immediately. Tissues, such as those from the

380 hypothalamus, pituitary, gonad, liver, kidney, spleen, rumen, reticulum, and small

381 intestine, were collected. RNA was extracted from the 15 tissue samples obtained from

382 the animals. After library construction and size selection, 150.4 Gb (57.7×) of long

383 reads with a mean length of 9,205 bp were generated by the PacBio RSII platform

384 (Single Molecule Real-Time, SMRT). In addition, 261.5 Gb (100.6×) of paired-end

385 data with varying insert sizes (200, 300, 400, and 600 bp) were generated by the

386 Illumina HiSeq 2000 platform (Additional file 1: Figure S17).

387 De novo *genome sequencing and Hi-C-based assembly*

388 The PacBio subreads were used to perform *de novo* genome assembly via wtdbg2 [17]

389 with the key parameter "-H –k 19". Then, primary assemblies were polished using the

390 Quiver [18] algorithm with the default parameters. A total of 93.4× clean paired-end

391 reads from the Illumina platform were aligned to the Quiver-polished assemblies using

392 BWA (v0.7.10-r943-dirty) to reduce the remaining InDel and base substitution errors

393 in the draft assembly. Inconsistent sequences between the polished genome and

394 Illumina reads were identified with SAMTools/VCFtools (v1.3.1). The credible

395 homozygous variations with differences in quality exceeding 20, a mapping quality

396 greater than 40 and a sum of high-quality alt-forward and alt-reverse bases more than

397 2 in the Quiver-polished assemblies were replaced by the called bases using in-house

398 scripts. Finally, highly accurate contigs were generated.

399 Four billion PE150 reads were produced from three Hi-C libraries by the Illumina

400 HiSeq platform. Hi-C-based proximity-guided scaffolding was used to connect primary

401 contigs. Clean reads were first aligned against the reference genome with the Bowtie2

402    end-to-end algorithm. HiC-Pro (v2.7.8) was then able to detect the ligation sites and

403    align them back to the genome with the 5' fraction of the reads. The assembly tool

404    LACHESIS was applied for clustering, ordering and orienting. Based on the

405    agglomerative hierarchical clustering algorithm, we clustered the contigs into 33 groups.

406    For each chromosome cluster, we obtained an exact scaffold order of the internal groups

407    and traversed all the directions of the scaffolds through a weighted directed acyclic

408    graph (WDAG) to predict the orientation for each scaffold. A chromosome-scale

409    assembly with 33 clusters was obtained that anchored 99.24% of the contigs for sika

410    deer.

411    *Genome accuracy assessment*

412    To determine the completeness and accuracy of the MHL_v1.0 assembly, we carried

413    out the following validation. First, the MHL_v1.0 assembly was aligned to the red deer

414    genome (CerEla1.0) and BioNano optical maps. The conflicting regions that appeared

415    in both alignments were potential misassemblies and were manually inspected

416    andcorrected.

417        A total of 2,715 EST sequences of sika deer were downloaded from the NCBI

418    dbEST database and aligned with MHL_v1.0 using BLAST (v35). The BUSCO [37]

419    software package was used to assess the quality of the generated genome using the

420    genome model "- M genome". The CEGMA pipeline software, which was also run

421    against the MHL_v1.0. Illumina short reads (93.4×), was aligned to MHL_v1.0 with

422    BWA to estimate the accuracy of a single base of the assembly, which was based on

423    the count of homozygous SNPs.

424    *Repeat sequence annotation*

425    To annotate the sika deer genome, RepeatModeler (v1.0.8) was initially used to obtain

426    a *de novo* repeat library. Next, RepeatMasker (v4.0.5) was used to search for known

427    and novel transposable elements (TEs) by mapping sequences against the Repbase TE

428    library (20150807) [38].

429    *Gene annotation*

430    For *de novo* gene prediction, we utilized AUGUSTUS (v3.0.3), SNAP (v2006-07-28),

431    GlimmerHMM (v3.0.4) and GENSCAN to analyze the repeat-masked genome. For

432    homology-based gene predictions, the protein sequences of human, mouse, cattle, sheep,

433    and horse were mapped to the sika deer genome with GenBlastA [39]. Then, prediction

434    was performed with GeneWise (v2.2.3) [40] in aligned regions. RNA-seq reads were

435    aligned to the genome using TopHat (v2.0.12) and assembled by Cufflinks (v 2.2.1)

436    with the default parameters. EVidenceModeler software (EVM, v1.1.1) was used to

437    integrate the genes predicted by homology, *de novo* and transcriptome approaches and

438    generate a consensus gene set. Short-length (< 50 aa) and transcriptome data for

439    nonsupport genes were removed from the consensus gene set, and the final gene set was

440    produced.

441    We translated the final predicted coding regions into protein sequences and mapped

442    all the predicted proteins to the Swiss-Prot, TrEMBL, and KEGG databases using

443    BLASTP (v2.2.27+) for gene functional annotation. We used the InterProScan database

444    to annotate the motifs, domains and Gene Ontology (GO) terms of proteins with

445    retrieval from the Pfam, PRINTS, PROSITE, ProDom, and SMART databases.

446    *Gene family construction*

447    Annotations of human, mouse, pig, sheep and cattle genomes were downloaded from

448    Ensembl (release-87), while those of minke whale, dromedary, Bactrian camel, yak,

449    goat, white-tailed deer, red deer, and reindeer were downloaded from NCBI. To

450    annotate the structures and functions of putative genes in the giraffe, okapi, milu, musk

451    deer, and roe deer assemblies, we used homology-based predictions. Cattle proteins

452    (Ensemble release-87) were aligned to the 5 genomes using GenBlastA (v1.0.1) [39]

453    and predicted by GeneWise (v2.2.3) [40]. The genes of the above 18 species and sika

454    deer were used to construct gene families using TreeFam [17]. All the protein sequences

455    were searched in the TreeFam (version 9) HMM file and classified among different

456    TreeFamilies.

457    *Phylogeny and divergence time estimation*

458    We constructed a phylogenetic tree based on a concatenated sequence alignment of 748

459    single-copy gene families from sika deer and 18 other mammalian taxa (human, mouse,

460    pig, sheep, cattle, minke whale, dromedary, Bactrian camel, yak, goat, white-tailed deer,

461    red deer, reindeer, giraffe, okapi, milu, musk deer, and roe deer) using the RAxML [41]

462    software with the GTRGAMMA model. Divergence times were estimated by PAML

463    [42] MCMCTREE. The Markov chain Monte Carlo (MCMC) process was run for

464    20,000 iterations with a sample frequency of 2 after a burn-in of 1,000 iterations. Other

465    parameters used the default settings of MCMCTREE. Two independent runs were

466    performed to check convergence. The following constraints were used for fossil time

467    calibrations: (1) Bovinae and Caprinae divergence time (18-22 Ma); (2) Ruminantia

468    and Suina divergence time (48.3-53.5 Ma); (3) Euarchontoglires and Laurasiatheria

469    divergence time (95.3-113 Ma); (4) Euarchontoglires and Rodentia divergence time

470    (85-94 Ma); and (5) Cervus and Elaphurus divergence time ($< 3$ Ma).

471    *Gene family expansions and contractions*

472    The CAFE program (v3.1) [43] was used to analyze gene family expansions and

473    contractions. The program uses a birth and death process to model gene gain and loss

474    across a user-specified phylogenetic tree. The numbers of sika deer genes relative to

475    the number of inferred ancestor genes and expanding and contracting gene families

476    were obtained. According to the GO and KEGG pathway results of the functional

477    annotation, the hypergeometric distribution was used for enrichment analysis, and the

478    BH (Benjamini and Hochberg) algorithm was used for *P*-value correction. A *P*-value

479    less than 0.05 after correction was considered a significant enrichment result.

480    We investigated several *UGT* genes in each category for the 19 species. The

481    annotated *UGT* genes of human and sika deer were used to predict the unannotated

482    *UGT* genes in the other 17 species with the program GeneWise [40]. MUSCLE

483    software was used for the multiple sequence alignment of all these *UGT* gene protein

484    sequences, whereby a phylogenetic *UGT* gene tree was constructed using RAxML [41].

485    *Synteny analysis*

486     A collinearity analysis between sika deer and red deer was conducted using the

487     MUMmer package (v3.23). Furthermore, to identify the synteny block among sika deer,

488     red deer, cattle and goats, we used MCscan (python version) [44] to search for and

489     visualize intragenomic syntenic regions. A homologous synteny block map between

490     sika deer and cattle was plotted with Circos.

491     *Demographic history reconstruction*

492     We inferred the demographic histories of sika deer using the Pairwise Sequentially

493     Markovian Coalescent (PSMC) model for diploid genome sequences. The whole-

494     genome diploid consensus sequence for PSMC input was generated by mapping short

495     reads to the sika deer genome with BWA (v0.7.10-r943-dirty) and SAMTools. Program

496     `fq2psmcfa' transforms the consensus sequence into a fasta-like format. The parameters

497     for `psmc' were set as follows: -N25 -t15 -r5 -p "4+25*2+4+6". The generation times

498     (g) of sika deer, cattle, and buffalo were 5 and 6 years, respectively. The mutation rate

499     for all species was 2.2e-9 per site per year.

500     *Positive selection genes*

501     For the single-copy orthologous genes of 19 species, multiple sequence alignment was

502     carried out using MUSCLE (v3.8.31). Regions of uncertain alignment were removed

503     by Gblocks 0.91b [45]. We used branch-site models and likelihood ratio tests (LRTs)

504     in the CODEML of PAML (v4.8a) [42] to detect positive selection genes (PSGs) in the

505     sika deer genome. *P*-values were computed using the $\chi^2$ statistic and corrected for

506     multiple testing by the false discovery rate (FDR) method (*Padj* < 0.05). All the PSGs

507     were mapped to KEGG pathways and assigned GO terms. GO and KEGG enrichment

508     analyses were then applied to detect the significantly enriched biological processes and

509     signaling pathways of positively selected genes (*Padj* < 0.05).

510     *Transcriptome analysis*

511     We performed RNA sequencing of 15 tissues (hypothalamus, liver, muscle, spleen,

512     kidney, testis, pituitary, cecum, duodenum, ileum, jejunum, rumen, abomasum,

513     reticulum and omasum) for each of the 12 sika deer from the feeding trials to determine

514 variations in gene expression levels after treatment. To compare the response to

515 different tannin levels between cattle and sika deer, we conducted RNA-seq and

516 transcriptome analyses of 8 tissues (hypothalamus, liver, kidney, rumen, jejunum,

517 pituitary, reticulum and spleen) from two groups of 6 individuals with a diet containing

518 0% or 10% gallotannic acid (GA). Total RNA from 226 feeding experiment samples

519 was extracted and used for library construction and sequencing. All libraries were

520 sequenced using an Illumina HiSeq platform.

521  The transcriptome data of each sample were mapped to the sika deer and cattle

522 genomes using HISAT2 (v2.0.5), and gene expression was calculated in each sample

523 using StringTie (v1.3.0). The R language package DESeq2 was used to homogenize the

524 expression and calculate the differential expression between each pair of samples, in

525 which genes with Padj < 0.05 were considered differentially expressed genes. For the

526 DEGs, the hypergeometric distribution and BH (Benjamin and Hochberg) algorithm

527 were used in the GO and KEGG enrichment analysis and *P*-value correction,

528 respectively. A Q value < 0.05 was considered significantly enriched in the GO and

529 KEGG pathways.

530

**Authors' contributions**

531  F.Y., X.X., C.L., and J.R. conceived of the project and designed the research; P.H.
532  drafted the manuscript with input from all authors; C.A., T.W., Y.L., H.T.L., Q.Q. and
533  Q. L. revised the manuscript; C.A., T.W., Y.L. and H.T.L. performed the majority of
534  the analysis, with contributions from H.M. L, R.Z., H.W. and L.W.; Y.C. and S.Z.
535  prepared the library and performed the sequencing; S.W. and A.L. performed the
536  genome assembly with help from W.Z.; T.W. obtained the Hi-C data; X.W. performed
537  the genome annotation analysis; C.A. conducted the positive selection and repeat
538  annotation analysis; X.S. and C.A. performed the gene family analysis; C.A. and T.W.
539  performed the genome collinearity analysis; T.W. performed the reverse transcription
540  analysis; G.W., H.T.L. and J.Z. conducted the feeding trials and prepared the samples
541  for transcriptome sequencing with help from H.W., R.Z., X.S., S.S., Z.Y., T.Y., Y.D.,
542  Y.J., L.S., P.Z., H.F., J.X. and X.C.; C.A., M.R., S.C., X.W., W.Y., M.Y.; T.W., and
543  Y.L. performed the analysis of transcriptome data; J.Y. and Y.T. provided the
544  geographic distribution data for Mongolian oak; Y.L. conducted the analysis of the
545  geographic distributions of Mongolian oak and sika deer; C.A., Y.L., T.W. and H.T.L.
546  performed the charting and graphing; and all authors read and approved the final
547  manuscript.

549

**Competing interests**

551  The authors declare no competing interests.

552

**Acknowledgements**

559 **Availability of data and material**

560 The whole-genome sequence data reported in this paper have been deposited in the

561 Genome Warehouse in the National Genomics Data Center, Beijing Institute of

562 Genomics (China National Center for Bioinformation), Chinese Academy of Sciences,

563 under accession number GWHANOY00000000, which is publicly accessible at

564 https://bigd.big.ac.cn/gwh. The raw sequence data have been deposited in the Genome

565 Sequence Archive in the National Genomics Data Center under accession numbers

566 CRA001393, CRA002054 and CRA002056, which are publicly accessible at

567 https://bigd.big.ac.cn/gsa.

568

569 All procedures concerning animals were performed in accordance with the guidelines

570 for the care and use of experimental animals established by the Ministry of Agriculture

571 of China, and all protocols were approved by the Institutional Animal Care and Use

572 Committee of Institute of Special Economic Animal and Plant Sciences, Chinese

573 Academy of Agricultural Sciences, Changchun, China.

574

575

# References

576

577 [1] Kierdorf U, Li C, Price JS. Improbable appendages: deer antler renewal as a
578 unique case of mammalian regeneration. Semin Cell Dev Biol 2009;20:535-42.
579 [2] Tseng SH, Sung CH, Chen LG, Lai YJ, Chang WS, Sung HC, et al. Comparison
580 of chemical compositions and osteoprotective effects of different sections of
581 velvet antler. J Ethnopharmacol 2014;151:352-60.
582 [3] Wu F, Li H, Jin L, Li X, Ma Y, You J, et al. Deer antler base as a traditional
583 Chinese medicine: a review of its traditional uses, chemistry and pharmacology.
584 J Ethnopharmacol 2013;145:403-15.
585 [4] Hillman JR, Davis RW, Abdelbaki YZ. Cyclic bone remodeling in deer. Calcif
586 Tissue Int 1973;12:323-30.
587 [5] Li C, Suttie JM, Clark DE. Morphological observation of antler regeneration in
588 red deer (*Cervus elaphus*). J Morphol 2004;262:731-40.
589 [6] Wang Y, Zhang C, Wang N, Li Z, Heller R, Liu R, et al. Genetic basis of
590 ruminant headgear and rapid antler regeneration. Science 2019;364:eaav6335.
591 [7] Doce RR, Hervás G, Belenguer A, Toral PG, Giráldez FJ, Frutos P. Effect of
592 the administration of young oak (*Quercus pyrenaica*) leaves to cattle on ruminal
593 fermentation. Anim Feed Sci Technol 2009;150:75-85.
594 [8] Li ZP, Liu HL, Li GY, Bao K, Wang KY, Xu C, et al. Molecular diversity of
595 rumen bacterial communities from tannin-rich and fiber-rich forage fed
596 domestic Sika deer (*Cervus nippon*) in China. BMC Microbiol 2013;13:151.
597 [9] Wan F, Yin C, Tang R, Chen M, Wu Q, Huang C, et al. A chromosome-level
598 genome assembly of *Cydia pomonella* provides insights into chemical ecology
599 and insecticide resistance. Nat Commun 2019;10:4237.
600 [10] Deschamps S, Zhang Y, Llaca V, Ye L, Sanyal A, King M, et al. A
601 chromosome-scale assembly of the *Sorghum* genome using nanopore
602 sequencing and optical mapping. Nat Commun 2018;9:4844.
603 [11] Johnson RN, O'Meally D, Chen Z, Etherington GJ, Ho SYW, Nash WJ, et al.
604 Adaptation and conservation insights from the koala genome. Nat Genet
605 2018;50:1102-11.
606 [12] Low WY, Tearle R, Bickhart DM, Rosen BD, Kingan SB, Swale T, et al.
607 Chromosome-level assembly of the water buffalo genome surpasses human and
608 goat genomes in sequence contiguity. Nat Commun 2019;10:260.
609 [13] Chen L, Qiu Q, Jiang Y, Wang K, Lin Z, Li Z, et al. Large-scale ruminant
610 genome sequencing provides insights into their evolution and distinct traits.
611 Science 2019;364:eaav6202.
612 [14] Lin Z, Chen L, Chen X, Zhong Y, Yang Y, Xia W, et al. Biological adaptations
613 in the *Arctic cervid*, the reindeer (*Rangifer tarandus*). Science
614 2019;364:eaav6312.
615 [15] Elsik CG, Tellam RL, Worley KC, Gibbs RA, Muzny DM, Weinstock GM, et
616 al. The genome sequence of taurine cattle: a window to ruminant biology and
617 evolution. Science 2009;324:522-8.

618  [16]  Zimin AV, Delcher AL, Florea L, Kelley DR, Schatz MC, Puiu D, et al. A
619        whole-genome assembly of the domestic cow, *Bos taurus*. Genome Biol
620        2009;10:R42.
621  [17]  Ruan J, Li H. Fast and accurate long-read assembly with wtdbg2. Nat Methods
622        2020;17:155-8.
623  [18]  Chin CS, Alexander DH, Marks P, Klammer AA, Drake J, Heiner C, et al.
624        Nonhybrid, finished microbial genome assemblies from long-read SMRT
625        sequencing data. Nat Methods 2013;10:563-9.
626  [19]  Bana NÁ, Nyiri A, Nagy J, Frank K, Nagy T, Stéger V, et al. kierThe red deer
627        *Cervus elaphus* genome CerEla1.0: sequencing, annotating, genes, and
628        chromosomes. Mol Genet Genom 2018;293:665-84.
629  [20]  Bickhart DM, Rosen BD, Koren S, Sayre BL, Hastie AR, Chan S, et al. Single-
630        molecule sequencing and chromatin conformation capture enable *de novo*
631        reference assembly of the domestic goat genome. Nat Genet 2017;49:643-50.
632  [21]  Venter JC, Adams MD, Myers EW, Li PW, Mural RJ, Sutton GG, et al. The
633        sequence of the human genome. Science 2001;291:1304-51.
634  [22]  Jiang Y, Xie M, Chen W, Talbot R, Maddox JF, Faraut T, et al. The sheep
635        genome illuminates biology of the rumen and lipid metabolism. Science
636        2014;344:1168-73.
637  [23]  Hu P, Shao Y, Xu J, Wang T, Li Y, Liu H, et al. Genome-wide study on genetic
638        diversity and phylogeny of five species in the genus *Cervus*. BMC Genom
639        2019;20:384.
640  [24]   Hedrick PW, McDonald JF. Regulatory gene adaptation: an evolutionary model.
641        Heredity 1980;45:83-97.
642  [25]  Li C, Harper A, Puddick J, Wang W, McMahon C. Proteomes and signalling
643        pathways of antler stem cells. PLoS One 2012;7:e30026.
644  [26]  Liu Z, Zhao H, Wang D, McMahon C, Li C. Differential effects of the
645        PI3K/AKT pathway on antler stem cells for generation and regeneration of
646        antlers *in vitro*. Front Biosci 2018;23:1848-63.
647  [27]  Meech R, Mackenzie PI. Structure and function of uridine diphosphate
648        glucuronosyltransferases. Clin Exp Pharmacol Physiol 1997;24:907-15.
649  [28]  Fedejko B, Mazerska Z. UDP-glucuronyltransferases in detoxification and
650        activation metabolism of endogenous compounds and xenobiotics. Postepy
651        Biochem 2011;57:49-62.
652  [29]  Wang H, Cao G, Wang G, Hao H. Regulation of mammalian UDP-
653        glucuronosyltransferases. Curr Drug Metab 2018;19:490-501.
654  [30]  Gilbert C, Ropiquet A, Hassanin A. Mitochondrial and nuclear phylogenies of
655        *Cervidae* (*Mammalia*, *Ruminantia*): systematics, morphology, and
656        biogeography. Mol Phylogenetics Evol 2006;40:101-17.
657  [31]  Feeny P, Bostock H. Seasonal changes in the tannin content of oak leaves.
658        Phytochemistry 1968;7:871-80.

659 [32] Li R, Fan W, Tian G, Zhu H, He L, Cai J, et al. The sequence and *de novo*
660    assembly of the giant panda genome. Nature 2010;463:311-7.

661 [33] Freedman AH, Gronau I, Schweizer RM, Ortega-Del Vecchyo D, Han E, Silva
662    PM, et al. Genome sequencing highlights the dynamic early history of dogs.
663    PLoS Genet 2014;10:e1004016.

664 [34] Kim JY, Cheong HS, Park BL, Kim LH, Namgoong S, Kim JO, et al.
665    Comprehensive variant screening of the UGT gene family. Yonsei Med J
666    2014;55:232-9.

667 [35] Doce R, Hervás G, Giráldez F, López-Campos O, Mantecón A, Frutos P. Effect
668    of immature oak (*Quercus pyrenaica*) leaves intake on ruminal fermentation
669    and adaptation of rumen microorganisms in cattle. J Anim Feed Sci 2007;16:13-
670    8.

671 [36] Kumar K, Chaudhary LC, Agarwal N, Kamra DN. Isolation and
672    characterization of tannin-degrading bacteria from the rumen of goats fed oak
673    (*Quercus semicarpifolia*) leaves. Agric Res 2014;3:377-85.

674 [37] Waterhouse RM, Seppey M, Simão FA, Manni M, Ioannidis P, Klioutchnikov
675    G, et al. BUSCO applications from quality assessments to gene prediction and
676    phylogenomics. Mol Biol Evol 2018;35:543-8.

677 [38] Bao W, Kojima KK, Kohany O. Repbase update, a database of repetitive
678    elements in eukaryotic genomes. Mob DNA 2015;6:11.

679 [39] She R, Chu JSC, Wang K, Pei J, Chen N. GenBlastA: enabling BLAST to
680    identify homologous gene sequences. Genome Res 2009;19:143-9.

681 [40] Birney E, Clamp M, Durbin R. Genewise and genomewise. Genome Res
682    2004;14:988-95.

683 [41] Stamatakis A. RAxML version 8: a tool for phylogenetic analysis and post-
684    analysis of large phylogenies. Bioinformatics 2014;30:1312-3.

685 [42] Yang Z. PAML: a program package for phylogenetic analysis by maximum
686    likelihood. Bioinformatics 1997;13:555-6.

687 [43] De Bie T, Cristianini N, Demuth JP, Hahn MW. CAFE: a computational tool
688    for the study of gene family evolution. Bioinformatics 2006;22:1269-71.

689 [44] Tang H, Wang X, Bowers JE, Ming R, Alam M, Paterson AH. Unraveling
690    ancient hexaploidy through multiply-aligned angiosperm gene maps. Genome
691    Res 2008;18:1944-54.

692 [45] Talavera G, Castresana J. Improvement of phylogenies after removing
693    divergent and ambiguously aligned blocks from protein sequence alignments.
694    Syst Biol 2007;56:564-77.

695

# Figure legends

**Figure 1    Distribution and genome assembly of sika deer**

**A**, Mongolian oak and sika deer distribution. The green shadow represents the distribution range of Mongolian oak. The yellow dots represent the historical distribution of sika deer in 5 countries (China, Russia, Japan, North Korea and Vietnam). **B**, A contact map at a 500-kb resolution of chromosome-level assembly in sika deer is shown. The color bar illuminates the logarithm of the contact density from red (high) to white (low) in the plot. Note that only sequences anchored on chromosomes are shown in the plot. **C**, Synteny analysis of cattle and sika deer. Circular graphs displaying the results from the synteny analysis. Same-color ribbons connect syntenic genomic segments.

**Figure 2    Evolutionary analysis of sika deer**

**A**, Phylogenetic tree inferred from 19 species. The x-axis is the inferred divergence time (M years) based on the phylogenetic tree and fossils. The number of expanded gene families is red, and the number of contracted gene families is blue. **B**, PSMC analysis of effective population sizes in sika deer, cattle and buffalo.

**Figure 3    *UGT* expansion and high-tannin adaptation in sika deer**

Transcriptome analysis revealed that the *UGT* gene family was the key factor for sika deer adaptation to a high-tannin diet. **A**, Gene tree of *UGTs* in 19 species. The red stars are significantly differentially expressed genes in the sika deer transcriptome. **B**, Number of *UGT* genes in 19 species. **C**, Expression heatmap of *UGTs* of sika deer in different tissues and treatments. **D**, The overlap between 3 contracted genes (yellow background), 20 expanded genes (green background) and 12 DEGs (pink background), which all play a role in the cytochrome P450 pathway. **E**, Expression change of 8 significant differentially expressed genes in sika deer liver resulting from different treatments. **F**, Six upregulated *UGT* genes in the *UGT2B* and *UGT2C* subfamilies were located on sika deer chromosome 27; schematic of the glucuronidation reaction. UDPGA, uridine 5'-diphospho-glucuronic acid.

724    **Tables**

725    **Table 1    Comparison of genome quality and annotation between the genome of**

726    **sika deer and the best published genome of red deer**

| | | Sika deer (Cervus nippon) | Red deer (Cervus elaphus) |
|---|---|---|---|
| Assembly | Total sequence length | 2,500,646,934 | 3,438,623,608 |
| | Total length without gaps | 2,500,501,634 | 1,960,832,178 |
| | Number of scaffolds | 588 | 11479 |
| | Scaffold N50/L50 | 78,786,809/12 | 107,358,006/13 |
| | Number of contigs | 2040 | 406637 |
| | Contig N50/L50 | 23,559,432/33 | 7,944/64532 |
| | Total number of chromosomes | 33 | 35 |
| | Anchored rate | 99.93% | 98.33% |
| Annotation | Gene number | 21499 | 19243 |
| | Average gene length | 39397.69 | 28008.84 |
| | Average CDS length | 1617.26 | 1085.04 |
| | Average exons per gene | 9.29 | 6.5 |
| | Average exon length | 174.03 | 167.06 |
| | Average intron length | 4555.82 | 4755.75 |

727

728 **Supplementary material**

729 **Figure S1   Karyotype of the sequenced female sika deer.** The karyotype analysis

730 shows that the sika deer chromosome number is 2n=66

731 **Figure S2   Distribution of the 25-mer frequency in the sika deer genome.** The

732 genome size of sika deer is 2.6 Gb based on Kmer analysis with Kmer=25

733 **Figure S3   Assembly strategy of the sika deer genome.** PacBio long reads were *de*

734 *novo* assembled with wtdbg2. The chromosome-scale scaffolds were generated by

735 using Hi-C data after genomic error correction. A BioNano optical map and proximal

736 species (red deer) genome were used to check the assembly accuracy

737 **Figure S4   Genome synteny analysis between sika deer and red deer.** The x-axis

738 represents red deer chromosomes, and the y-axis represents sika deer chromosomes.

739 These two assemblies show significant genomic synteny

740 **Figure S5   Hi-C interaction heatmap for each chromosome of the sika deer**

741 **genome**

742 **Figure S6   Gene syntenic blocks between the sika deer genome and the three**

743 **ruminant genomes.** The representative chromosome fission/separation fragment is

744 indicated in purple, turquoise and cyan. Gray wedges in the background highlight

745 conserved syntenic blocks with more than 10 gene pairs

746 **Figure S7   Distribution of identified transposable elements among different**

747 **mammalian species.** Data anomalies of red deer may be due to the poor quality of the

748 genome

749 **Figure S8   Circos plot of the chromosomal features of sika deer.** The external

750 green circle represents the chromosomes of sika deer. The circles and links inside the

751 chromosomes from outside to inside represent the distribution of genes in the

752 chromosomes (blue); distribution of repeats of the genome (orange); distribution of

753 heterozygosity (green); and segmental duplications (length >10 kb) (red)

754 **Figure S9   Phylogeny and divergence time of 19 species.** Maximum-likelihood (ML)

755 tree inferred from single-copy orthologous genes by RAxML. The x-axis is the inferred

756 divergence time (M year) based on the phylogenetic tree and fossils

757 **Figure S10   Gene family expansion and contraction analysis.** The number of

758 expanded gene families is in red, and the number of contracted gene families is in green

759 **Figure S11   Phylogenetic tree of all *UGT* genes.** Phylogeny structured by RAxML

760 based on the multiple sequence alignment of all *UGT* genes. These *UGTs* were divided

761 into seven groups. The star represents significantly differentially expressed genes

762 **Figure S12   Expression heatmap of differentially expressed genes (DEGs) among**

763 **different treatments**

764 **Figure S13   Expression of *UGT* genes in 8 tissues of cattle.** *UGT* genes were highly

765 expressed in the liver, kidney and jejunum

766 **Figure S14   *CYP* gene expression patterns in sika deer.** Five differentially

767 expressed *CYP* genes were upregulated in the liver tissue with increasing tannin intake

768 **Figure S15   Potential metabolism of drugs and exogenous substances, such as**

769 **tannins, in the mammalian body.** Oak leaves are rich in hydrolysable tannins. Proline-

770 rich salivary proteins (PRPs) found in the mouth can precipitate gallotannic acid (GA)

771 and play a role in the defense against GA. However, PRPs are not found in all the

772 published genomes of cattle, sheep and our Mhl_v1.0. In the rumen, GA is hydrolyzed

773 into gallic acid and ellagic acid, which are degraded by rumen microbes into simple

774 phenolic compounds. Some of these compounds can be metabolized by the P450

775 enzyme and excreted from the body. Glucuronyltransferase (GT), sulfatyltransferase

776 (SULT), glutathione S-transferase (GST) and other enzymes produced by the liver can

777 catalyze the conversion of undigested phenolic compounds into glucuronates, sulfates

778 and other water-soluble compounds that can be excreted through the urine. Our results

779 show that only the expression of *UGTs* increased with the tannin content in the liver

780 **Figure S16   Comparison of the liver, kidney and heart in sika deer, cattle and**

781 **sheep after a tannin feeding experiment.** The three tissues showed no difference

782    between the treatment group and the control group in sika deer. However, lesions (white

783    arrow) occurred in the three tissues of cattle and sheep. These results demonstrated

784    different tannin tolerances among the 3 species

785    **Figure S17   Distribution of the insertion segment of Illumina paired-end data**.

786    Illumina sequencing data were generated with four different insert fragment sizes (200,

787    300, 400, and 600 bp)

788    **Table S1   Estimation of the sika deer genome size using K-mer analysis**

789    **Table S2   Summary of the genome sequencing of sika deer**

790    **Table S3   Summary of the sika deer genome assembly**

791    **Table S4   Summary of the Hi-C assembly of chromosome length in sika deer**

792    **Table S5   Summary of the Cervidae genome assembly**

793    **Table S6   Assessment of the completeness and accuracy of the sika deer genome**

794    **Table S7   Summary of the repeat content in the sika deer genome**

795    **Table S8   Comparison of the identified transposable elements among different**

796    **mammalian species**

797    **Table S9   Functional annotation of sika deer genes**

798    **Table S10   Summary of predicted protein-coding genes and gene characteristics**

799    **Table S11   BUSCO of annotation and assembly**

800    **Table S12   Statistics for the gene families**

801    **Table S13   Positively selected genes (PSGs) identified in sika deer**

802    **Table S14   Functionally enriched KEGG pathway categories of sika deer**

803    **expanded genes**

804    **Table S15   Functionally enriched GO categories of sika deer expanded genes**

805    **Table S16   Functionally enriched KEGG pathway categories of sika deer**
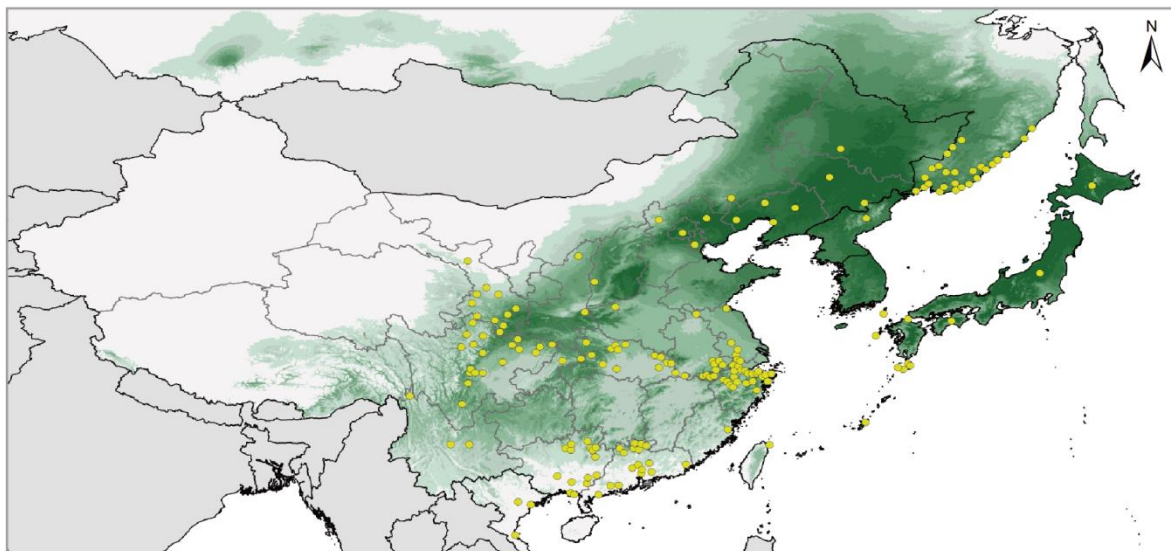
806    **contracted genes**

807    **Table S17   Functionally enriched GO categories of sika deer contracted genes**

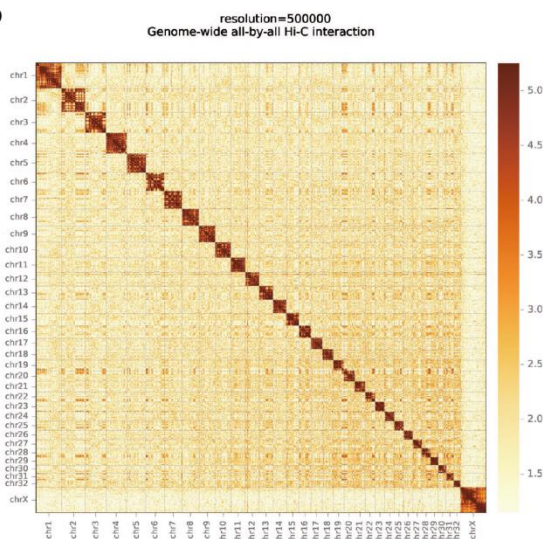808    **Table S18   Numbers of annotated *UGT* genes in 19 species**

809    **Table S19   Design of the feeding experiment**
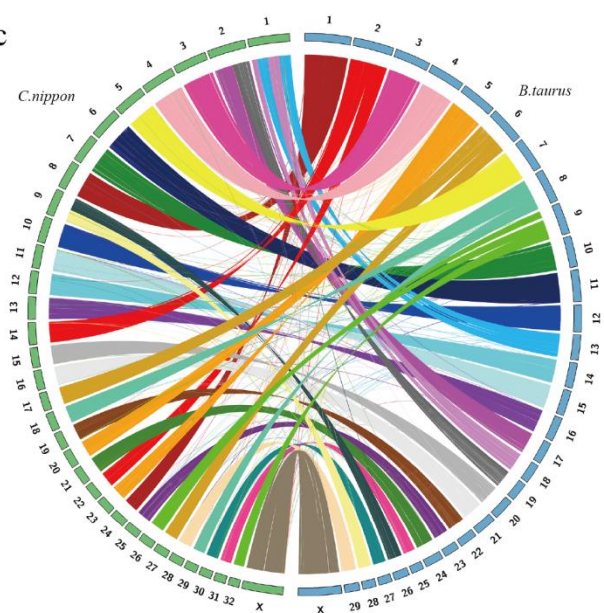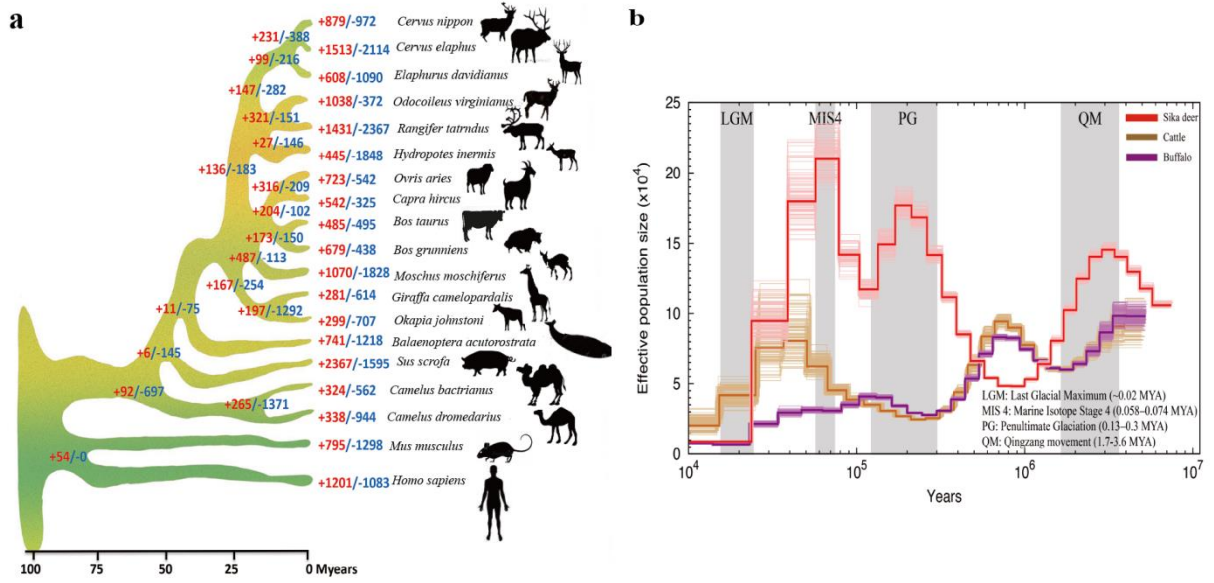
810 **Figures**

a



b

resolution=500000
Genome-wide all-by-all Hi-C interaction
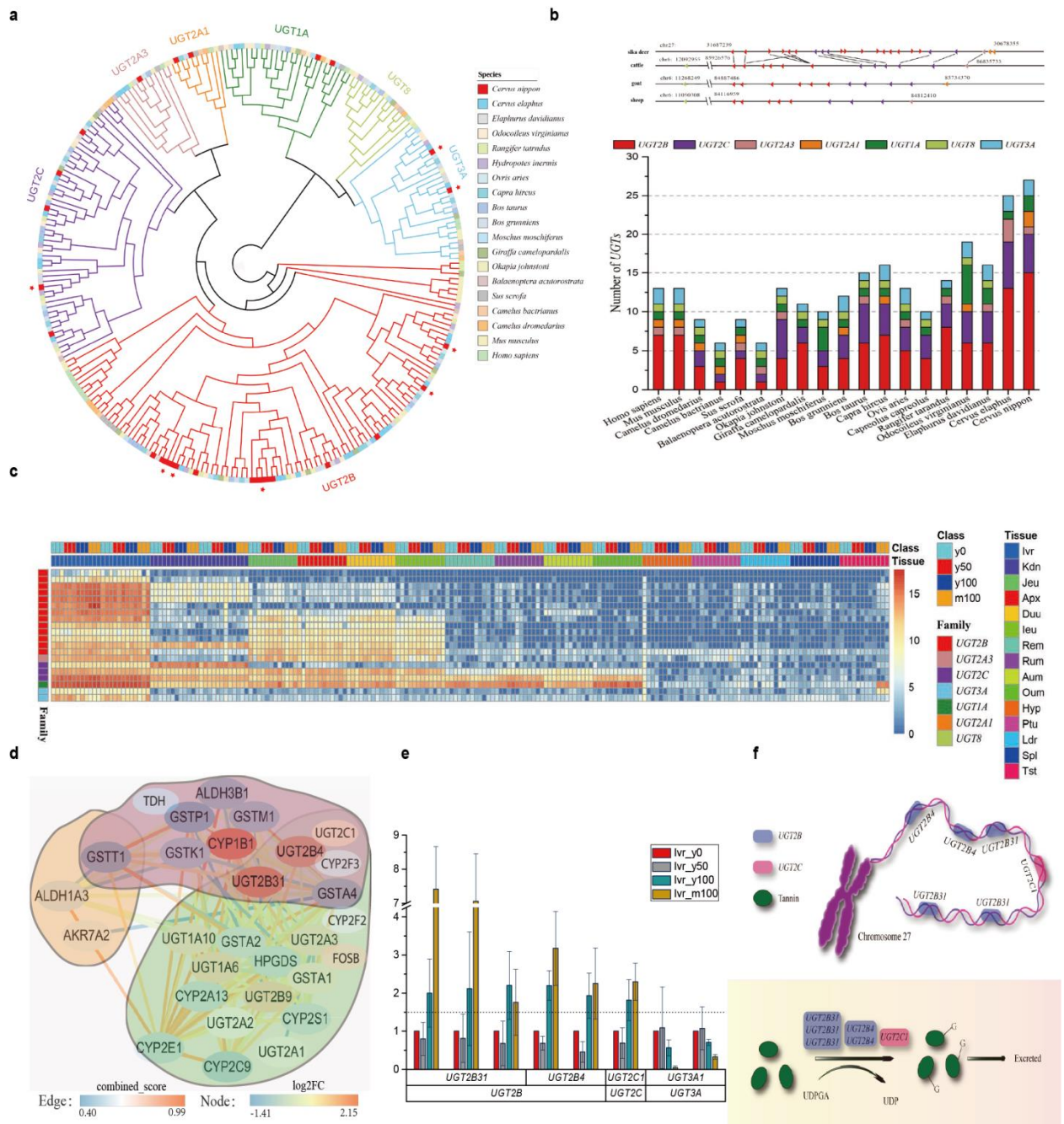
c

C.nippon        B.taurus

811

812 **Figure 1**

813

814

**Figure 2**

816

817

818 **Figure 3**