**DATABASE**                                                                                    **Open Access**

# Yak genome database: a multi-omics analysis platform

Hui Jiang[1,2†], Zhi-Xin Chai[3†], Xiao-Ying Chen[1,2†], Cheng-Fu Zhang[1,2], Yong Zhu[1,2], Qiu-Mei Ji[1,2*] and Jin-Wei Xin[1,2*]

## Abstract

**Background**  The yak (*Bos grunniens*) is a large ruminant species that lives in high-altitude regions and exhibits excellent adaptation to the plateau environments. To further understand the genetic characteristics and adaptive mechanisms of yak, we have developed a multi-omics database of yak including genome, transcriptome, proteome, and DNA methylation data.

**Description**  The Yak Genome Database (http://yakgenomics.com/) integrates the research results of genome, transcriptome, proteome, and DNA methylation, and provides an integrated platform for researchers to share and exchange omics data. The database contains 26,518 genes, 62 transcriptomes, 144,309 proteome spectra, and 22,478 methylation sites of yak. The genome module provides access to yak genome sequences, gene annotations and variant information. The transcriptome module offers transcriptome data from various tissues of yak and cattle strains at different developmental stages. The proteome module presents protein profiles from diverse yak organs. Additionally, the DNA methylation module shows the DNA methylation information at each base of the whole genome. Functions of data downloading and browsing, functional gene exploration, and experimental practice were available for the database.

**Conclusion**  This comprehensive database provides a valuable resource for further investigations on development, molecular mechanisms underlying high-altitude adaptation, and molecular breeding of yak.

**Keywords**  Yak, Genome, Database, Multi-omics, Plateau environment

†Hui Jiang, Zhi-Xin Chai and Xiao-Ying Chen contributed equally to this work.

*Correspondence:
Qiu-Mei Ji
jiqiumei07@163.com
Jin-Wei Xin
xinjinwei18@163.com
¹State Key Laboratory of Hulless Barley and Yak Germplasm Resources and Genetic Improvement, 850000 Lhasa, Tibet, China
²Institute of Animal Science and Veterinary, Tibet Academy of Agricultural and Animal Husbandry Sciences, 850000 Lhasa, Tibet, China
³Key Laboratory of Qinghai-Tibetan Plateau Animal Genetic Resource Reservation and Utilization, Sichuan Province and Ministry of Education, Southwest Minzu University, 610041 Chengdu, Sichuan, China

## Background

Although single omics study provides information and insights into specific biological or molecular processes, it is hard to confirm the real molecular mechanisms underlying the functionality of an organism and the relationships between biological processes and environmental factors. Integrating and analyzing multiple omics data provide an effective and systematic approach to life science researchers. In general, genomics provides DNA sequence information, transcriptomics examines gene transcription patterns under specific conditions, proteomics explores the composition and expression levels of proteins in cells, and DNA methylation involves chemical modifications on DNA molecules [1]. Multi-omics

Jiang *et al. BMC Genomics*      (2024) 25:346

Page 2 of 8

analysis combines data at different levels to comprehensively explore biological processes. Multi-omics analysis reveals connections between genomics, transcriptomics, proteomics, and DNA methylation data, facilitating to understand how genomic variations impact gene transcription and protein expression, as well as the associations between DNA methylation and gene activities [2]. These pieces of information contribute novel information to the gene regulatory networks, which are important to molecular mechanisms underlying biological functions, development, metabolism, etiopathology, and environmental adaptation.

The yak (*Bos grunniens*) is a unique species in the Qinghai-Tibet Plateau, and widely distributes in high-altitude areas of Western China and neighboring regions. As a large mammal at the highest-altitude area, yak has survived and adapted to the harsh and cold environment after thousands of years of evolution [3]. Their unique biological features make them an ideal model for studying adaptive evolution and high-altitude ecosystems. Yak also plays important roles in agriculture and economic development. As a significant livestock species, yak provides meat, fur, and other economic resources. Their dung is also an important source of agricultural fertilizer and energy production. Moreover, yak positively impacts the ecological balance and vegetation restoration in the plateau grasslands through their grazing behaviors [4]. In recent years, we have analyzed yak using different omics approaches. These data preliminarily explored the yak genetic characteristics, gene transcription, protein expression, and DNA methylation patterns, as well as molecular regulatory mechanisms in response to different conditions [5–11], providing novel insights into the mechanisms underlying evolution, and high-altitude adaptation in yak.

Currently, the data resources of yak omics researches are generally stored in public databases in their raw data format, such as NCBI. These databases primarily provide storage and retrieval functions, but lack an integrated platform for data integration and in-depth analysis. Hu et al. [12] developed a yak genome database (http://me.lzu.edu.cn/yak), which incorporated genome sequences, predicted genes and associated annotations, non-coding RNA sequences, transposable elements, and single nucleotide variants of yak, as well as three-way whole-genome alignments between human, cattle and yak. However, this database did not include other omics datasets, such as transcriptome, proteome, and DNA methylation. Given the vast and diverse nature of omics data, the traditional database retrieval methods could not fully explore the relationship between different types of datasets [13]. Thus, an integrated platform of different omics data is crucial to facilitate data integration, interaction, and analysis. An integrated platform can also offer advanced data mining and machine learning algorithms to help researchers discover the complex relationships among yak genomics, transcriptomics, proteomics, and other omics levels, further deepening our understanding of biological processes and diseases in yak.

In this study, the Yak Genome Database (http://yak-genomics.com/) was constructed, which successfully assembled a comprehensive yak fine-scale genome map at the chromosome level, using PacBio sequencing, Illumina sequencing, Bionano assembly, and Hi-C three-dimensional genome scaffolding. Moreover, this platform also integrated transcriptome, proteome, and DNA methylation data of yak, which were not available in Yak Genome Database developed by Hu et al. [12]. This database provides basic information for yak researches in future, such as molecular breeding, molecular evolution, disease prevention and control.

## Construction and content

The Yak Genome Database was deployed in the Ubuntu 20.04 operation system using the AKKA 2.13 (web server), MySQL 8.0.30 (database server), Scala 2.13.2, and SBT 1.3.9. All data were managed and stored using the MySQL Database Management System. The query function was enforced based on Slick 3.3.2 middleware tier. The Jbrowse 1.16.11 was used to visualize the genome. The website interfaces were designed and implemented using the Bootstrap 4.6.0 and the Play Framework 2.8.7. The software versions and statistical tools used for data analyses and plot preparation have been presented in Xin et al. [6–11]. The boxplots, and heatmaps were prepared using R 4.2.1. The website has been tested in several popular web browsers, including Firefox, Google Chrome, and Internet Explorer.

## Utility and discussion

### The yak genome database content

The multi-omics data in the Yak Genome Database are categorized into two central functional domains: data resources and navigation (Fig. 1). The data resources contain four main modules, including genome, transcriptome, proteome, and methylation information. The database contains 26,518 genes, 62 transcriptomes, 144,309 proteome spectra, and 22,478 methylation sites of yak. The navigation page consists of Browser, Jbrowse, Search and Blast functions. Currently, the database supports individual download of images and gene data. In the future, we will add functions such as one-click download of whole genome information.

### Genome module

The Genome module incorporates the complete genomic DNA sequence of yaks obtained by the third-generation high-throughput sequencing platform (PacBio RSII) [14].
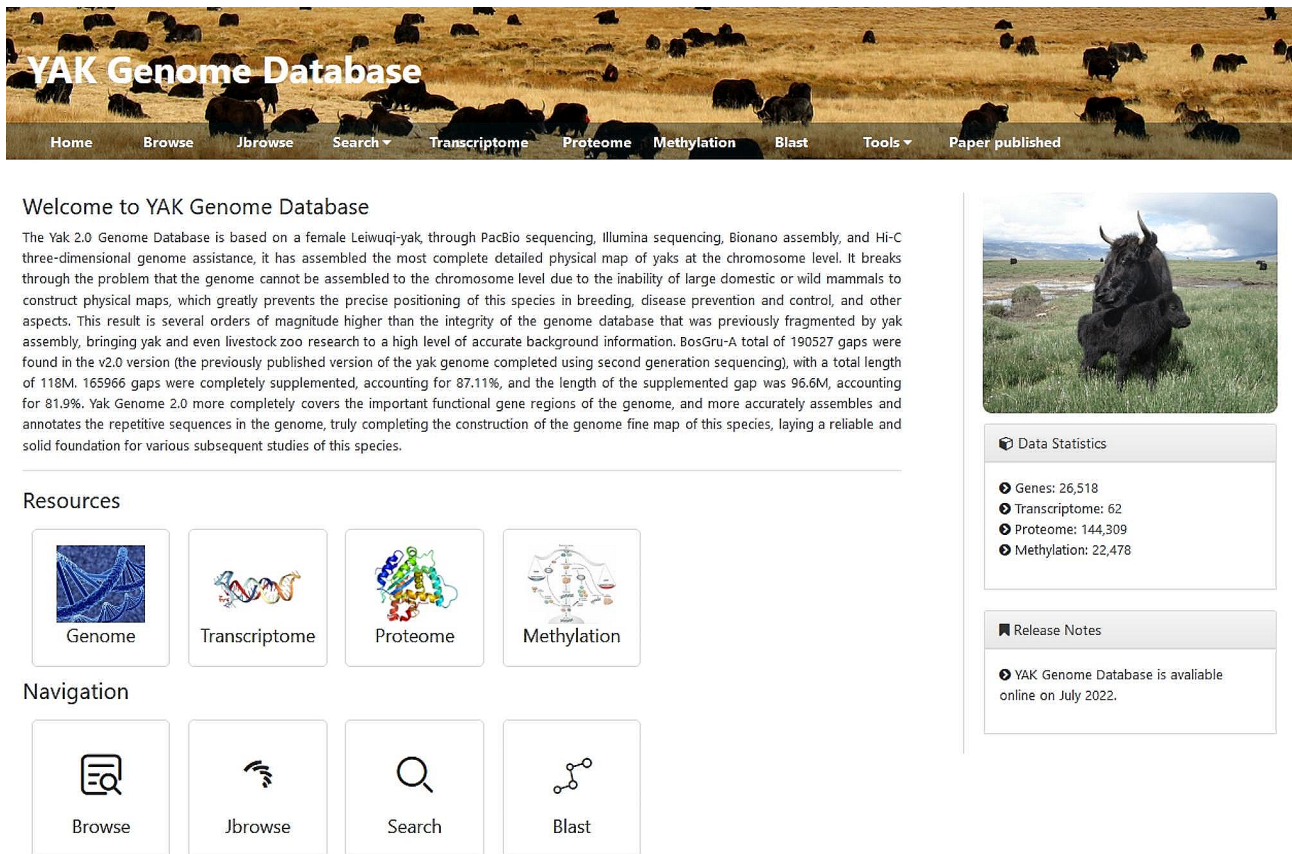
**Fig. 1** The homepage of yak genome database

The yak genome was sequenced at a coverage of 70X, with the second-generation sequencing data used to correct errors. The Bionano assisted assembly technology was used for high-quality assembly, and analysis. Next, a refined physical map of the yak chromosome was generated, providing a more readable and complete genome database than the fragmented information in another Yak Genome Database (BosGru_v2.0) [15], and contributing a novel genome tool to yak researchers.

When accessing the 'Genome' section on the homepage, a new page will display information of genes at all locations, such as Gene ID, Chromosome, Start Position, End Position, Strand, GO (Gene Ontology) terms, Interpro, KEGG (Kyoto Encyclopedia of genes and Genomes), Swissprot, and Trembl in a user-friendly table format (Fig. 2A). When clicking each gene, users can access detailed information of this gene, including annotations, transcriptional levels, proteome data, Jbrowse page, and nucleotide sequences associated with the gene (Fig. 2B-2D). The 'Annotation' tav provides comprehensive gene annotation information, including GO terms, KEGG pathways, and Interpro annotations, which can be further explored by clicking them. The 'Expression' tab displays gene expression levels across different cattle breeds and tissues, and users can download the images in various formats by selecting the menu in the upper right corner of the image. 'Jbrowse' is used to display integrated information from annotated genomic datasets, while 'Seqs' provides the coding sequence (CDS) and protein sequence on the selected gene.

**Transcriptome module**

Previously, comparative transcriptome sequencing was performed on lung, gluteal muscle, and mammary gland tissues of low-altitude cattle (Sanjiang and Holstein cattle), Tibetan cattle (living at a moderate altitude), and yaks (living at a high altitude). In addition, these tissues of yaks at different ages (6, 30, 60, and 90 months) were also subjected to transcriptome sequencing. These analyses identified the functional genes involved in the major biochemical, metabolic, and signal transduction pathways involved in yak development and high-altitude adaptation [10, 11]. These data are included in the transcriptome module on the website, providing a valuable transcriptome database for specific tissue biomarkers, molecular research, and breeding of yaks. After clicking the "Transcriptome" button, users can select the strain in the 'Sample' dialog box, enter the gene ID in the 'Gene ID' dialog box, and then click 'Search' (Fig. 3A), and then the website will return the transcriptional levels of the
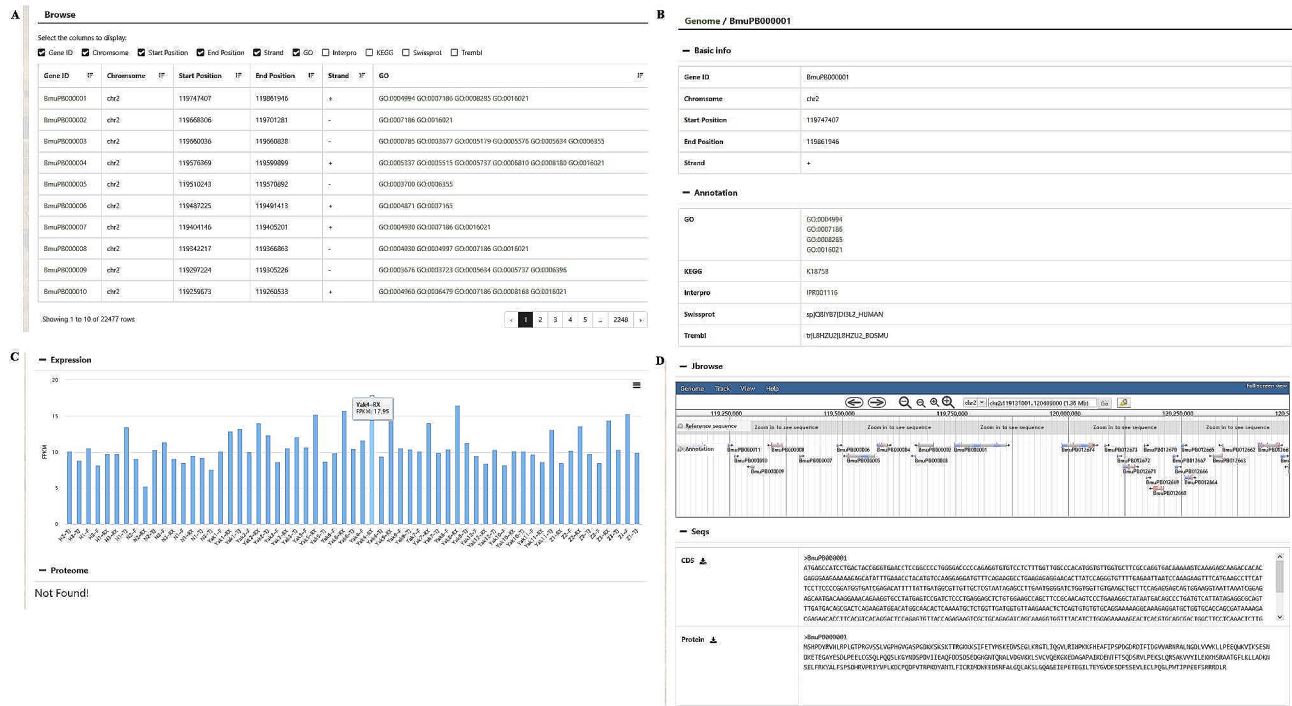
**Fig. 2** Features of the genome module. (**A**) Genome browse. (**B**) Basic information and annotation of a gene. (**C**) Gene expression. (**D**) Gene Jbrowse and sequences
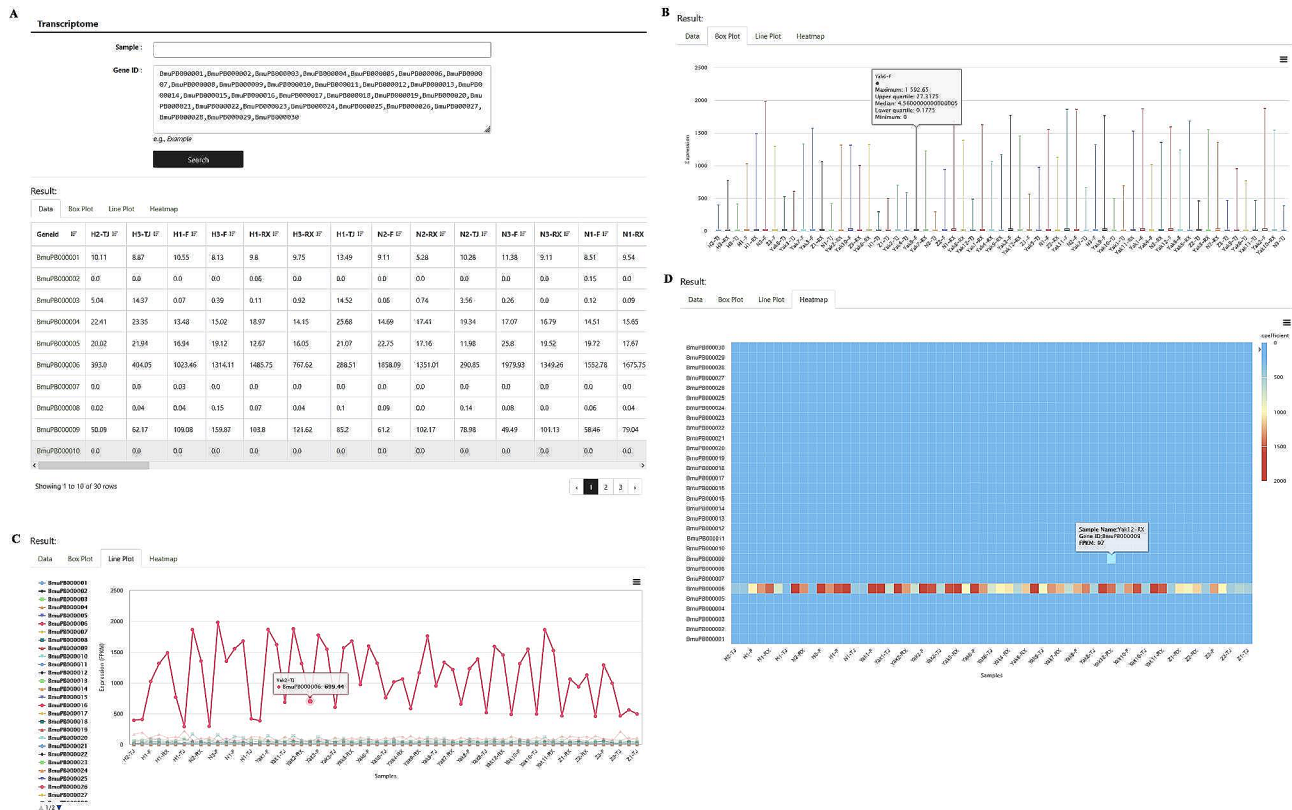


**Fig. 3** Features of the transcriptome module. (**A**) Transcriptome browse. (**B**) Box plot, (**C**) Line plot and (**D**) Heatmap of gene expression

selected genes in selected samples in the forms of data table, Boxplot, Lineplot, and Heatmap (Fig. 3B and D).

## Proteome module

Using the liquid chromatography-mass spectrometry (LC-MS) method, proteomic analyses were conducted for four specific tissues from four different species (yak, Tibetan cattle, Sanjiang cattle, and Holstein cattle) [7–9]. All the animals were female and 60 months of age. The proteome module provides two input dialog boxes. Users can select two samples and then click the "search" button. Next, the website will return the comparison results of the expression levels of all genes in the two selected samples, including log2(fold change) and statistical parameters (Fig. 4).

## Methylation module

DNA methylation is a critical epigenetic modification that occurs in both animals and plants, playing pivotal roles in chromosome structure, gene expression and regulation [16]. The establishment of a comprehensive DNA methylation database for yak can significantly advance the comprehension of cellular gene expression and regulation, and provide deeper insights into the spatiotemporal specificity of DNA methylation across various developmental stages and organs [17]. The DNA

methylation database of yak presents single-base methylation maps and tissue-specific methylation maps. The single-base methylation maps include: 1) DNA methylation levels at the single-base resolution, 2) DNA methylation levels specific to different base types, 3) DNA methylation levels specific to different gene structures, 4) DNA methylation levels in repetitive sequences, and 5) DNA methylation levels in non-coding sequences and regulatory regions. The tissue-specific methylation maps involve three tissues: mammary gland, lung, and muscle [6]. On the website, users can select 'Sample' and 'Chromosome' in the Methylation module, set the 'Start Position' and 'End Position,' and finally click 'Search' to obtain the corresponding DNA methylation results on the selected sequences (Fig. 5).

## Navigation

'Browse' allows users to read the yak genome directly. 'JBrowse' is a next-generation genome browser built with JavaScript and HTML5. The Jbrowse of Yak Genome Database includes tracks describing gene, gene sequence, mRNAs, structure, and other gene-related features, and provides a graphical display of annotations on the yak genome (Fig. 6). Users can browse gene models on chromosomes and unanchored contigs. For example, if user set the genomic region from 4,454,001 bp to 5,878,000 bp
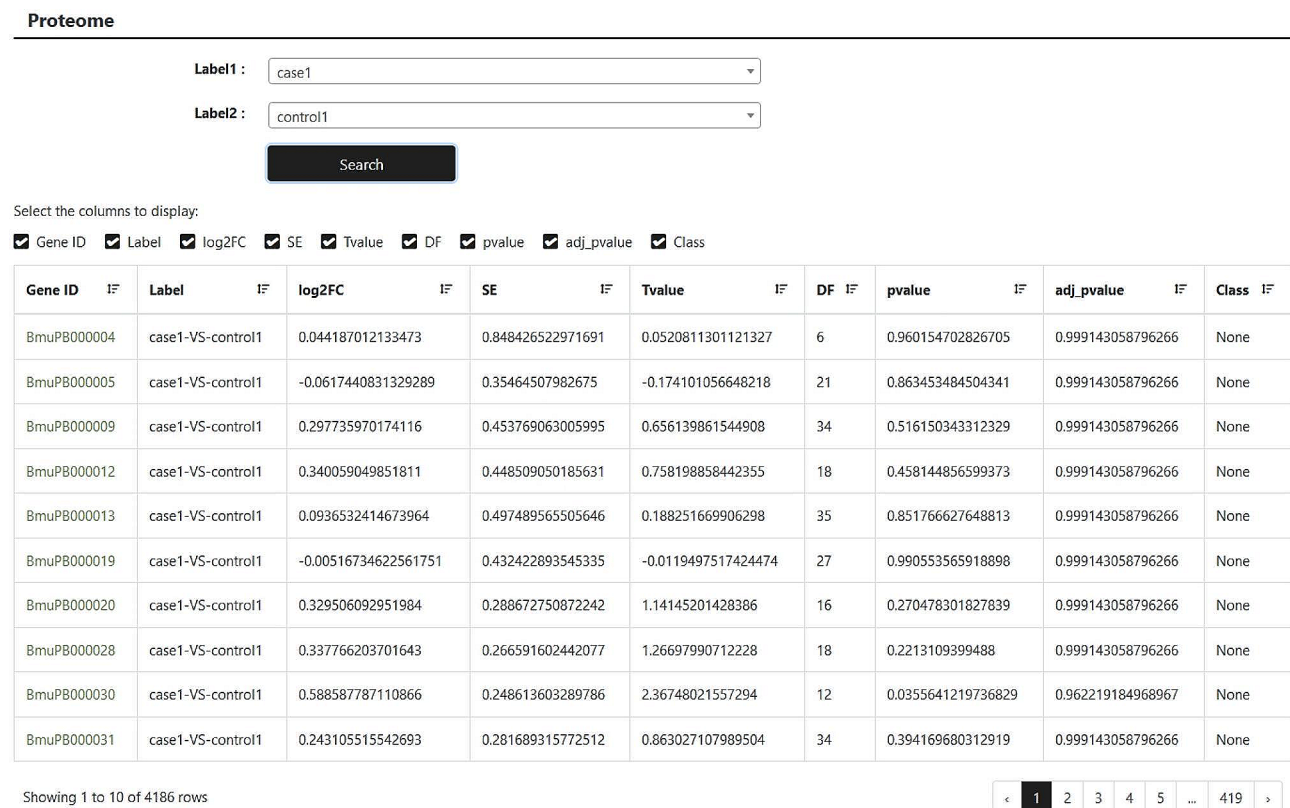
---

**Proteome**

**Label1 :**   [ case1                                                          ▾ ]

**Label2 :**   [ control1                                                       ▾ ]

[ **Search** ]

Select the columns to display:

☑ Gene ID  ☑ Label  ☑ log2FC  ☑ SE  ☑ Tvalue  ☑ DF  ☑ pvalue  ☑ adj_pvalue  ☑ Class

| Gene ID ⇅ | Label ⇅ | log2FC ⇅ | SE ⇅ | Tvalue ⇅ | DF ⇅ | pvalue ⇅ | adj_pvalue ⇅ | Class ⇅ |
|---|---|---|---|---|---|---|---|---|
| BmuPB000004 | case1-VS-control1 | 0.044187012133473 | 0.848426522971691 | 0.0520811301121327 | 6 | 0.960154702826705 | 0.999143058796266 | None |
| BmuPB000005 | case1-VS-control1 | -0.0617440831329289 | 0.35464507982675 | -0.174101056648218 | 21 | 0.863453484504341 | 0.999143058796266 | None |
| BmuPB000009 | case1-VS-control1 | 0.297735970174116 | 0.453769063005995 | 0.656139861544908 | 34 | 0.516150343312329 | 0.999143058796266 | None |
| BmuPB000012 | case1-VS-control1 | 0.340059049851811 | 0.448509050185631 | 0.758198858442355 | 18 | 0.458144856599373 | 0.999143058796266 | None |
| BmuPB000013 | case1-VS-control1 | 0.0936532414673964 | 0.497489565505646 | 0.188251669906298 | 35 | 0.851766627648813 | 0.999143058796266 | None |
| BmuPB000019 | case1-VS-control1 | -0.00516734622561751 | 0.432422893545335 | -0.0119497517424474 | 27 | 0.990553565918898 | 0.999143058796266 | None |
| BmuPB000020 | case1-VS-control1 | 0.329506092951984 | 0.288672750872242 | 1.14145201428386 | 16 | 0.270478301827839 | 0.999143058796266 | None |
| BmuPB000028 | case1-VS-control1 | 0.337766203701643 | 0.266591602442077 | 1.26697990712228 | 18 | 0.2213109399488 | 0.999143058796266 | None |
| BmuPB000030 | case1-VS-control1 | 0.588587787110866 | 0.248613603289786 | 2.36748021557294 | 12 | 0.0355641219736829 | 0.962219184968967 | None |
| BmuPB000031 | case1-VS-control1 | 0.243105515542693 | 0.281689315772512 | 0.863027107989504 | 34 | 0.394169680312919 | 0.999143058796266 | None |

Showing 1 to 10 of 4186 rows                                      ‹  **1**  2  3  4  5  …  419  ›

**Fig. 4** Browse of the proteome module

**Fig. 5** Features of the methylation module. (**B**) Box plot, (**C**) Line plot and (**D**) Heatmap of methylated gene expression

on Chr1 for browsing, all genes in this region will appear in order (Fig. 6A). When clicking on 'BmuPB021145', an extra layer will appear with the detailed information, such as mRNAs, CDS and other features (Fig. 6B). For more operational details, users can click the 'Help' button, which provides comprehensive instructions and guidance.

The 'Search' tab supplies users with two methods (search by gene ID or range) for genome searching. When users click on 'Blast', three options 'Blastn Gene', 'Blastn Genome' and 'Blastp' will display. Users can select the Blast type and enter a DNA or protein sequence, and set the parameters of 'Evalue', 'Word size' and 'Max target seqs'. After clicking the 'Search' button, the nucleotide or protein sequence complying the search conditions will display and could be downloaded by the users.

### Additional tools
The Yak Genome Database also provides users with several convenient online tools, including Primer designer, GO and KEGG enrichment. The 'Primer designer' tool offers primer design function to amplify a selected sequence. The 'GO enrichment' and 'KEGG enrichment' tools facilitate the users to obtain the GO and KEGG enrichment results of a set of genes.

### Maintenance of the yak genome database in future
To ensure continuous operation of the Yak Genome Database, we would assign an administrator to manage the website regularly. We would keep omics studies on yak in future, and all the omics data we obtained would be uploaded to this database. In addition, we would keep cooperations with other investigators and find more cooperators who work on yak. Next, all the progresses on yak omics would also be encouraged to supplement in this database.

### Conclusions
The Yak Genome Database is a comprehensive platform of genomic physical map, which integrates genome, transcriptome, proteome, and DNA methylation data. Information in the database can be downloaded, and shared through the Internet. Users who want to upload their own data can contact the administrator of the website. By providing timely updates on yak research progress, the Yak Genome Database enables efficient and interactive sharing of existing scientific data among researchers worldwide who are interested in yak, cattle, livestock, ruminant animals, and even medical research. Comparative analysis of multidimensional data from key yak tissues aims to uncover the mechanisms underlying high-altitude adaptation, disease resistance, cold tolerance, and starvation resistance of large animals in the plateau. These findings contribute to molecular breeding
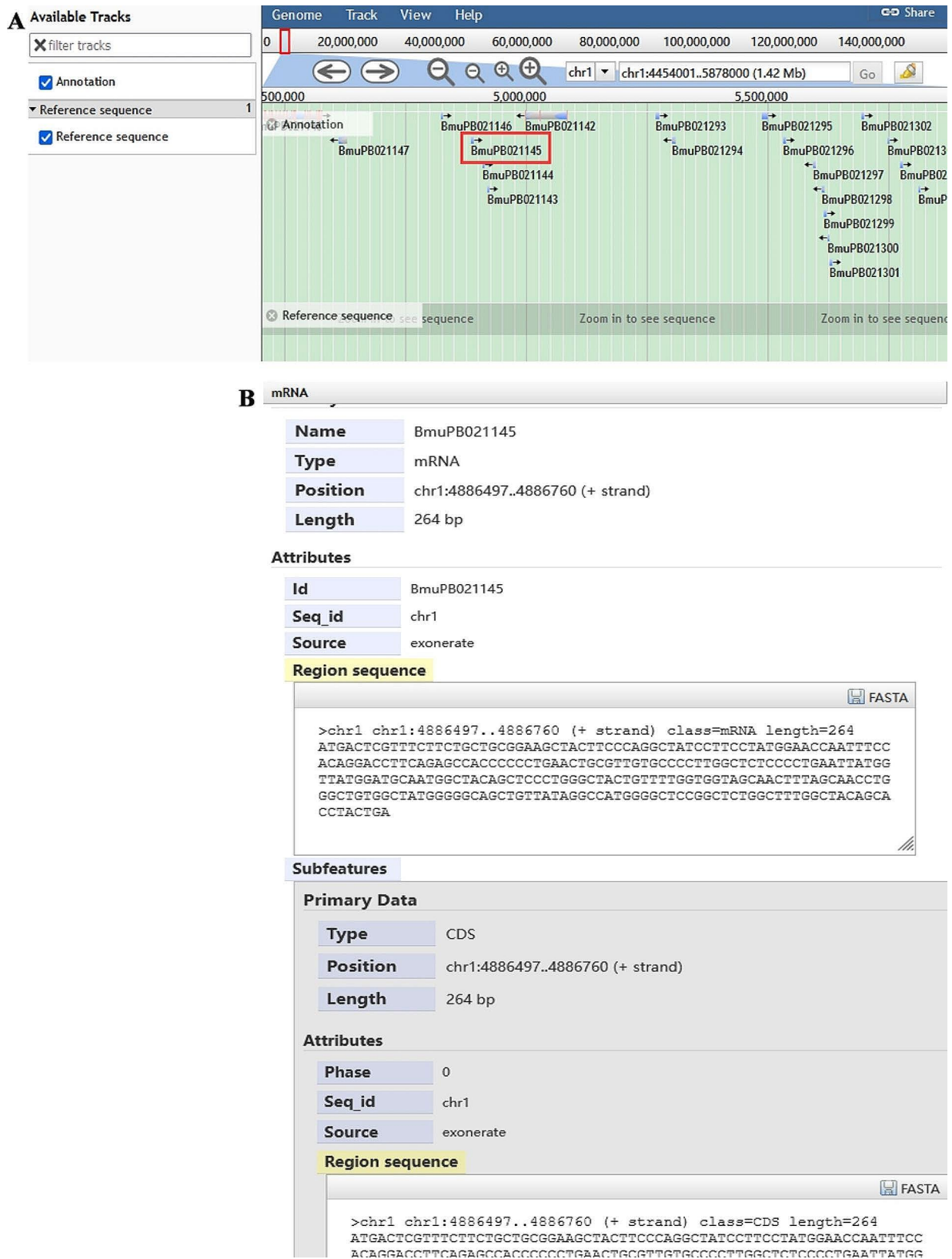
**Fig. 6** Regional view of the genome using Jbrowse. (**A**) A graphic view of the region 4,454,001 bp to 5,878,000 bp on Chr1. (**B**) The interface after clicking on 'BmuPB021145'.

Jiang *et al. BMC Genomics*　　　(2024) 25:346

Page 8 of 8

of livestock animals and the understanding of human responses to harsh environments.

## Abbreviations

| | |
|---|---|
| CDS | Coding Sequence |
| GO | Gene Ontology |
| KEGG | Kyoto Encyclopedia of Genes and Genomes |
| NCBI | National Center for Biotechnology Information |

## Data availability
The datasets generated and analyzed in the current study are freely available on the Download page of Yak database with the web link: http://yakgenomics.com/.

## Declarations

### Ethics approval and consent to participate
All procedures and experiments involving animals followed the guidelines for the Care and Use of Laboratory Animals. The Ethics Committee at Institute of Animal Science and Veterinary, Tibet Academy of Agricultural and Animal Husbandry Sciences (Permit Number: 2015 − 216) approved this study.

### Consent for publication
Not applicable.

### Competing interests
The authors declare no competing interests.

## References
1. Manzoni C, Kia DA, Vandrovcova J, Hardy J, Wood NW, Lewis PA, et al. Genome, transcriptome and proteome: the rise of omics data and their integration in biomedical sciences. Brief Bioinform. 2018;19(2):286–302.
2. Liao Y, Wang J, Zou J, Liu Y, Liu Z, Huang Z. Multi-omics analysis reveals genomic, clinical and immunological features of SARS-CoV-2 virus target genes in pan-cancer. Front Immunol. 2023;14:1112704.
3. Ge Q, Guo Y, Zheng W, Zhao S, Cai Y, Qi X. Molecular mechanisms detected in yak lung tissue via transcriptome-wide analysis provide insights into adaptation to high altitudes. Sci Rep. 2021;11(1):7786.
4. Ayalew W, Chu M, Liang C, Wu X, Yan P. Adaptation mechanisms of Yak (*Bos grunniens*) to high-Altitude Environmental stress. Animals. 2021;11(8):2344.
5. Gao X, Wang S, Wang YF, Li S, Wu SX, Yan RG, et al. Long read genome assemblies complemented by single cell RNA-sequencing reveal genetic and cellular mechanisms underlying the adaptive evolution of yak. Nat Commun. 2022;13(1):4887.
6. Xin J, Chai Z, Zhang C, Zhang Q, Zhu Y, Cao H, et al. Methylome and transcriptome profiles in three yak tissues revealed that DNA methylation and the transcription factor ZGPAT co-regulate milk production. BMC Genom. 2020;21(1):731.
7. Xin JW, Chai ZX, Zhang CF, Zhang Q, Zhu Y, Cao HW, et al. Signature of high altitude adaptation in the gluteus proteome of the yak. J Exp Zool B Mol Dev Evol. 2020;334(6):362–72.
8. Xin JW, Chai ZX, Zhang CF, Zhang Q, Zhu Y, Cao HW, et al. Differences in proteomic profiles between yak and three cattle strains provide insights into molecular mechanisms underlying high-altitude adaptation. J Anim Phys Anim Nutr. 2022;106(3):485–93.
9. Xin JW, Chai ZX, Zhang CF, Yang YM, Zhang Q, Zhu Y, et al. Comparative analysis of Skeleton muscle Proteome Profile between Yak and cattle provides insight into high-altitude adaptation. Curr Proteom. 2021;18(1):62–70.
10. Xin JW, Chai ZX, Zhang CF, Zhang Q, Zhu Y, Cao HW, et al. Transcriptome profiles revealed the mechanisms underlying the adaptation of yak to high-altitude environments. Sci Rep. 2019;9(1):7558.
11. Xin JW, Chai ZX, Zhang CF, Zhang Q, Zhu Y, Cao HW, et al. Comparisons of lung and gluteus transcriptome profiles between yaks at different ages. Sci Rep. 2019;9(1):14213.
12. Hu Q, Ma T, Wang K, Xu T, Liu J, Qiu Q. The yak genome database: an integrative database for studying yak biology and high-altitude adaption. BMC Genomics. 2012;13:600.
13. Tarazona S, Arzalluz-Luque A, Conesa A. Undisclosed, unmet and neglected challenges in multi-omics studies. Nat Comp Sci. 2021;1(6):395–402.
14. Ji QM, Xin JW, Chai ZX, Zhang CF, Dawa Y, Luo S, et al. A chromosome-scale reference genome and genome-wide genetic variations elucidate adaptation in yak. Mol Ecol Res. 2021;21(1):201–11.
15. Jiangfeng F, Yuzhu L, Sijiu Y, Yan C, Gengquan X, Libin W, et al. Transcriptional profiling of two different physiological states of the yak mammary gland using RNA sequencing. PLoS ONE. 2018;13(7):e0201628.
16. Lucibelli F, Valoroso MC, Aceto S, Plant DNA, Methylation. An epigenetic Mark in Development, Environmental interactions, and evolution. Int j mol sci. 2022;23(15):8299.
17. Chai Z, Wu Z, Ji Q, Wang J, Wang J, Wang H, et al. Genome-wide DNA methylation and hydroxymethylation changes revealed epigenetic regulation of Neuromodulation and Myelination in Yak Hypothalamus. Front Genet. 2021;12:592135.

## Publisher's Note
Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.