

GigaScience

Comparative transcriptomics of five high-altitude vertebrates and their low-altitude relatives

--Manuscript Draft--

Manuscript Number:	GIGA-D-17-00037R2	
Full Title:	Comparative transcriptomics of five high-altitude vertebrates and their low-altitude relatives	
Article Type:	Data Note	
Funding Information:	National High Technology Research and Development Program of China (863 Program) (2013AA102502)	Prof. Mingzhou Li
	National Natural Science Foundation of China (31402046)	Dr Qianzi Tang
	National Natural Science Foundation of China (31522055)	Prof. Mingzhou Li
	National Natural Science Foundation of China (31601918)	Dr Jideng Ma
	National Natural Science Foundation of China (31530073)	Prof Xuewei Li
	National Natural Science Foundation of China (31472081)	Prof. Mingzhou Li
	Science & Technology Support Program of Sichuan (2016NYZ0042)	Dr Yiren Gu
	Youth Science Fund of Sichuan (2017JQ0011)	Dr Yiren Gu
	China Postdoctoral Science Foundation (2015M572486)	Dr Qianzi Tang
	China Agriculture Research System (CARS-36)	Dr Yiren Gu
	Program for Innovative Research Team of Sichuan Province (2015TD0012)	Prof. Mingzhou Li
	Program for Pig Industry Technology System Innovation Team of Sichuan Province (SCCXTD-005)	Dr Yiren Gu
	Project of Sichuan Education Department (15ZA0008)	Dr Xun Wang
	Project of Sichuan Education Department (15ZA0003)	Dr Miaomiao Mai
	Project of Sichuan Education Department (16ZA0025)	Dr Jideng Ma
	Project of Sichuan Education Department (16ZB0037)	Dr An'an Jiang
	National Program for Support of Top-notch Young Professionals	Prof. Mingzhou Li
	Young Scholars of the Yangtze River	Prof. Mingzhou Li
	National Natural Science Foundation of China (31772576)	Prof. Mingzhou Li
Abstract:	Background: Species living at high altitude are subject to strong selective pressures due to inhospitable environments (e.g., hypoxia, low temperature, high solar radiation, and lack of biological production), making these species valuable models for comparative analyses of local adaptation. Studies that examined high-altitude	

	<p>adaptation identified a vast array of rapidly evolving genes that characterize the dramatic phenotypic changes in high-altitude animals. However, how high-altitude environment shapes gene expression programs remains largely unknown.</p> <p>Findings: We generated a total of 910 Gb high-quality RNA-seq data for 180 samples derived from six tissues of five agriculturally important high-altitude vertebrates (Tibetan chicken, Tibetan pig, Tibetan sheep, Tibetan goat and yak), and their cross-fertile relatives living in geographically neighboring low-altitude regions. Of these, ~75% reads could be aligned to their respective reference genomes, and on average ~70% of annotated protein coding genes in each organism showed FPKM expression values greater than 0.1. We observed a general concordance in topological relationships between the nucleotide alignments and gene expression-based trees. Tissue and species accounted for markedly more variance than altitude based on either the expression or the alternative splicing patterns. Cross-species clustering analyses showed a tissue-dominated pattern of gene expression, and a species-dominated pattern for alternative splicing. We also identified numerous differentially expressed genes were potentially involved in phenotypic divergence shaped by high-altitude adaptation.</p> <p>Conclusions: This data serves as a valuable resource for examining the convergence and divergence of gene expression changes between species as they adapt or acclimatize to high-altitude environments.</p> <p>Keywords: high-altitude vertebrates, comparative transcriptomics, gene expression, alternative splicing</p>
Corresponding Author:	Mingzhou Li, Ph.D. Sichuan Agricultural University Chengdu, Sichuan CHINA
Corresponding Author Secondary Information:	
Corresponding Author's Institution:	Sichuan Agricultural University
Corresponding Author's Secondary Institution:	
First Author:	Qianzi Tang
First Author Secondary Information:	
Order of Authors:	Qianzi Tang
	Yiren Gu
	Xuming Zhou
	Long Jin
	Jiuqiang Guan
	Rui Liu
	Jing Li
	Keren Long
	Shilin Tian
	Tiandong Che
	Silu Hu
	Yan Liang
	Xuemei Yang
	Xuan Tao
	Zhijun Zhong
	Guosong Wang
	Xiaohui Chen

	Diyang Li
	Jideng Ma
	Xun Wang
	Miaomiao Mai
	An'an Jiang
	Xiaolin Luo
	Xuebin Lv
	Vadim N. Gladyshev
	Xuwei Li
	Mingzhou Li, Ph.D.
Order of Authors Secondary Information:	
Response to Reviewers:	<p>Reviewer 1:</p> <p>Comment 2-1 Your Series GSE93855 submission is labeled as unavailable until Jan 13 2020. Under the FAIR principles of publication in GigaScience, this must change if paper is accepted. Other datasets are also embargoed, far into 2018 or 2019.</p> <p>Response 2-1 Thank you for your reminder. We have released all of datasets related to our manuscript, inducing a total of 180 RNA-seq data deposited in the NCBI Gene Expression Omnibus (GEO) under accession numbers GSE93855, GSE77020 and GSE66242, as well as the 7 whole-genome sequencing data deposited in the NCBI-sequence read archive (SRA) under accession number SRP096151.</p> <p>GSE93855: https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE93855 GSE77020: https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE77020 GSE66242: https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE66242 SRP096151: https://trace.ncbi.nlm.nih.gov/Traces/sra/?study=SRP096151</p> <p>Comment 2-2 The additional work performed to answer my reviewer comments is very good and shows interesting results. I don't have further major concerns, other than fixing the data availability issue above.</p> <p>Response 2-2 Thank you for your positive comments.</p> <p>Comment 2-3 The authors could consider revising Figure 4b and 4d, which are quite complex and difficult to see the relationships of all datasets to each other. The use of different color lines was not explained, and I tried to see why these were used, but was unable to see the logic.</p> <p>Response 2-3 Thank you for your valuable suggestion. In Fig. 4b and 4d, the vertical leading lines with different colors from the plotted points dropping to the xy-plane serve to improve the demonstration of group separation by tissue based on the first and second principal components. We added the explanation of the dropping vertical lines to the legends of Fig. 4b and 4d, "The vertical leading lines with different colors from the plotted points dropping to the xy-plane show the separation of points based on the first and second principal components". To further promote the clarity of depiction for the principal component analysis (PCA) results, we added the two-dimensional PCA figures as the Supplementary Figs. S14-15 (accessible from Response_sup.pdf at: https://www.dropbox.com/s/sivc1djpqvqzty2s/Response_sup.pdf?dl=0)</p> <p>Reviewer 2:</p> <p>Comment 3-1</p>

	<p>Thank you for addressing the previous points raised during the first reviewing stage. I believe that most points have been addressed in an appropriate way, and once again I thank you for your work and the data generated.</p> <p>Response 3-1 Thank you for your positive comments.</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
<p>Resources</p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p>	Yes
<p>Availability of data and materials</p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the "Availability of Data and Materials" section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our Minimum</p>	Yes

[Standards Reporting Checklist?](#)

1 Comparative transcriptomics of five high-altitude 2 vertebrates and their low-altitude relatives

3 Qianzi Tang^{1†}, Yiren Gu^{2†*}, Xuming Zhou^{3†}, Long Jin¹, Jiuqiang Guan⁴, Rui Liu¹, Jing Li¹,
4 Kereng Long¹, Shilin Tian¹, Tiandong Che¹, Silu Hu¹, Yan Liang², Xuemei Yang², Xuan
5 Tao², Zhijun Zhong², Guosong Wang^{1,5}, Xiaohui Chen², Diyan Li¹, Jideng Ma¹, Xun
6 Wang¹, Miaomiao Mai¹, An'an Jiang¹, Xiaolin Luo⁴, Xuebin Lv², Vadim N. Gladyshev³,
7 Xuewei Li¹ and Mingzhou Li^{1*}

8 ¹ Institute of Animal Genetics and Breeding, College of Animal Science and Technology,
9 Sichuan Agricultural University, Chengdu 611130, China;

10 ² Animal Breeding and Genetics Key Laboratory of Sichuan Province, Pig Science Institute,
11 Sichuan Animal Science Academy, Chengdu 610066, China

12 ³ Division of Genetics, Department of Medicine, Brigham and Women's Hospital, Harvard
13 Medical School, Boston, Massachusetts, 02115 USA;

14 ⁴ Yak Research Institute, Sichuan Academy of Grassland Science, Chengdu 610097,
15 China;

16 ⁵ Department of Animal Science, Texas A & M University, College Station, Texas, 77843
17 USA.

18 † These authors contributed equally to this work.

19 Corresponding authors. E-mail: Mingzhou Li: mingzhou.li@sicau.edu.cn, Yiren Gu:
20 guyiren1128@163.com.

21

22

23

1 24

2
3 25

4
5
6 26

Abstract

7
8 27 **Background:** Species living at high altitude are subject to strong selective
9
10 28 pressures due to inhospitable environments (e.g., hypoxia, low temperature,
11
12 29 high solar radiation, and lack of biological production), making these species
13
14 30 valuable models for comparative analyses of local adaptation. Studies that
15
16 31 examined high-altitude adaptation identified a vast array of rapidly evolving
17
18 32 genes that characterize the dramatic phenotypic changes in high-altitude
19
20 33 animals. However, how high-altitude environment shapes gene expression
21
22 34 programs remains largely unknown.

23
24
25
26
27
28 35 **Findings:** We generated a total of 910 Gb high-quality RNA-seq data for 180
29
30 36 samples derived from six tissues of five agriculturally important high-altitude
31
32 37 vertebrates (Tibetan chicken, Tibetan pig, Tibetan sheep, Tibetan goat and yak),
33
34 38 and their cross-fertile relatives living in geographically neighboring low-altitude
35
36 39 regions. Of these, ~75% reads could be aligned to their respective reference
37
38 40 genomes, and on average ~60% of annotated protein coding genes in each
39
40 41 organism showed FPKM expression values greater than 0.5. We observed a
41
42 42 general concordance in topological relationships between the nucleotide
43
44 43 alignments and gene expression-based trees. Tissue and species accounted
45
46 44 for markedly more variance than altitude based on either the expression or the
47
48 45 alternative splicing patterns. Cross-species clustering analyses showed a
49
50 46 tissue-dominated pattern of gene expression, and a species-dominated pattern
51
52 47 for alternative splicing. We also identified numerous differentially expressed
53
54 48 genes that could potentially be involved in phenotypic divergence shaped by
55
56
57
58
59
60
61
62
63
64
65

49 high-altitude adaptation.

50 **Conclusions:** This data serves as a valuable resource for examining the
51 convergence and divergence of gene expression changes between species as
52 they adapt or acclimatize to high-altitude environments.

53 **Keywords:** high-altitude vertebrates, comparative transcriptomics, gene
54 expression, alternative splicing

55

56 **Data description**

57 ***Transcriptome sequencing***

58 Six tissues (heart, kidney, liver, lung, skeletal muscle and spleen) of three
59 unrelated adult females for each of five high-altitude vertebrates and their low-
60 altitude relatives were sampled (**Fig. 1a** and **Supplementary Fig. S1**). Animals
61 were sacrificed humanely to ameliorate suffering. All animals and samples used
62 in this study were collected according to the guidelines for the care and use of
63 experimental animals established by the Ministry of Agriculture of China. We
64 extracted total RNA, prepared libraries and sequenced the libraries on Illumina
65 HiSeq 2000 or 2500 platforms. We generated a total of ~909.6 Gb high-quality
66 RNA-seq data for 180 samples (~5.05 Gb per sample) of 30 individuals across
67 6 tissues (**Supplementary Table S1**).

68 ***Whole-genome re-sequencing***

69 To compare the phylogeny derived from gene expression with the
70 phylogenetic relationships of the five high-altitude vertebrates and their low-
71 altitude relatives, we constructed the phylogenetic tree based on nucleotide
72 alignments. We extracted the unassembled reads from short-insert (500 bp)
73 libraries of a single yak [1] (NCBI-SRA: SRX103159 to SRX103161, and

74 SRX103175 and SRX103176), a Tibetan pig [2] (NCBI-SRA: SRX219342) and
75 a low-altitude Rongchang pig (NCBI-SRA: SRX1544519) [3] that were used for
76 *de novo* assemblies to roughly 10 × depth coverage. We also randomly
77 selected an individual of the cattle, low- and high-altitude chicken, goat and
78 sheep, and sequenced their whole genomes at ~10 × depth coverage (NCBI-
79 SRA: SRP096151). Genomic DNA was extracted from blood tissue of each
80 individual. Sequencing was performed on the Illumina X Ten platform, and a
81 total of 198.64 Gb of paired-end DNA sequence was generated
82 **(Supplementary Table S2).**

83

84 **Data analysis**

85 ***Data filtering***

86 To avoid reads with artificial bias, we removed the following type of reads: (a)
87 Reads with ≥ 10% unidentified nucleotides (N); (b) Reads with > 10 nt aligned
88 to the adapter, allowing ≤ 10% mismatches; (c) Reads with > 50% bases having
89 phred quality < 5.

90 ***Identification of single-copy orthologous genes***

91 Single-copy orthologous genes across five reference genomes, i.e. chicken
92 (Galgal4) [4], pig (Suscrofa 10.2) [5], cattle (UMD3.1) [6], goat (CHIR_1.0) [7]
93 and sheep (Oar_v3.1) [8] were determined using a EnsemblCompara
94 GeneTrees method [9] **(Supplementary Fig. S2, Supplementary Methods)**
95 **[9].**

96 ***Construction of phylogenetic tree based on nucleotide alignments***

97 High-quality re-sequencing data were mapped to their respective reference

98 genomes using BWA software, version 0.7.7 (BWA, RRID:SCR_010910) [10],
99 reads with mapping quality > 0 were retained and potential PCR duplication
100 cases were removed. For each individual, ~97.01% of reads were mapped to
101 ~97.40% (at least 1 × depth coverage) or ~91.86% (at least 4 × depth coverage)
102 of the reference genome assemblies (**Supplementary Table S2**).
103 Single nucleotide variations (SNVs) and insertion and deletions (InDels) were
104 further detected by following GATK's best practice, version 3.3-0 (GATK,
105 RRID:SCR_001876) [11]. We substituted SNVs and InDels identified in our
106 study in the coding DNA sequences (CDS) of the respective reference
107 genomes. Single copy orthologues with substituted CDS of the five vertebrates
108 were applied to Treebest [12] and generating the neighbor-joining tree (Fig. 1b).

109 *Analyses of gene expression*

110 High-quality RNA-seq reads were mapped to their respective reference
111 genomes using Tophat version 2.0.11 (TopHat, RRID:SCR_013035) [13].
112 Cufflinks version 2.2.1 (Cufflinks, RRID:SCR_014597) [14] was applied to
113 quantify gene expression and obtain FPKM expression values. We generated
114 abundance files by applying Cuffquant (part of Cufflinks) to read mapping
115 results. Log₂-transformed values of (FPKM + 1) for genes with >0.5 FPKM in
116 over 80% of the samples were used for subsequent analyses.

117 Pearson's correlations were calculated across six samples from low- and
118 high-altitudes populations within each group of specific tissue and animals;
119 among pairwise comparisons of five animals within each of the six tissues; and
120 among pairwise comparisons of six tissues within each of the five animals.
121 Principal Variance Component Analysis (PVCA) was carried out using R
122 package pvca [15]. Neighbor-joining expression-based trees were generated

123 according to distance matrices composed of pairwise (1-Spearman's
124 correlations) implemented in the R package ape [16]. Reproducibility of
125 branching patterns was estimated by bootstrapping genes, that is, the single
126 copy orthologues were randomly sampled with replacement 100 times. The
127 fractions of replicate trees that share the branching patterns of the original tree
128 constructed were marked by distinct node colors in the figure.

129 We generated abundance files by applying Cuffquant (part of Cufflinks) to
130 read mapping results, and further applied abundance files to Cuffdiff (part of
131 Cufflinks) to detect DEGs between population pairs from distinct altitudes
132 within each group of specific tissue and species. Genes with FDR-adjusted p-
133 values ≤ 0.05 were detected as DEGs.

134 Genes were converted to human orthologs, and assessed by DAVID
135 (DAVID, RRID:SCR_001881) [17] webserver for functional enrichment in GO
136 (Gene Ontology) terms consisting of molecular function (MF) and biological
137 process (BP) as well as the KEGG (KEGG, RRID:SCR_012773) pathways
138 and InterPro (InterPro, RRID:SCR_006695) databases (Benjamini adjusted p-
139 values ≤ 0.05).

140 ***Analyses of alternative splicing***

141 Single-copy orthologous exons were identified by finding annotated exons that
142 overlapped with the query exonic region in a multiple alignment of 99 vertebrate
143 genomes including human genome (hg38) from the UCSC genome browser
144 [18]. Exon groups with multiple overlapping exons in any species were
145 excluded. Each internal exon in every annotated transcript was taken as an

146 “cassette” exon. Each “cassett” alternative splicing (AS) is composed of three
147 exons: C1, A and C2, where A is alternative exon, C1 the 5’ alternative exon,
148 C2 the 3’ alternative exon. For each species and read length k, we generated
149 all non-redundant constitutive and alternative junction sequences for the
150 following RNA-seq alignments. The junction sequences were constructed by
151 retrieving k-8 bp from each of the two exons making up the junction, and when
152 the exon length is smaller than k-8, the whole sequence of the exon is retrieved.
153 This ensures that there is at least 8 bp overlap between the mapped reads and
154 each of the two junction exons.

155 We then estimated the effective number of uniquely mappable positions of
156 the junctions. We extracted L-k+1 (L being the junction length) k-mers from
157 each junction and mapped such k-mers back to the reference genome allowing
158 up to two mismatches. Those k-mers that failed to align were further mapped
159 to the non-redundant junctions. The number of k-mers that could uniquely align
160 to a junction was counted and deemed as the effective number of uniquely
161 mappable positions for the junction.

162 For each sample, RNA-seq reads were first aligned to the reference genome
163 allowing up to two mismatches, and the unaligned reads were further mapped
164 to the non-redundant junctions. Uniquely mapped reads for each junction were
165 counted, and multiplied by the ratio between the maximum number of mappable
166 positions (i.e. k-15) to the effective number of uniquely mappable positions for
167 the junction.

168 The “percent-spliced in” (PSI) values for each internal exon was defined as
169 $PSI = 100 \times \text{average}(\#C1A, \#AC2) / (\#C1C2 + \text{average}(\#C1A, \#AC2))$, here
170 #C1A, #AC2 and #C1C2 are the normalized read counts for the associated

171 junctions. Exons were taken as alternative in a sample if $5 \leq \text{PSI} \leq 95$. We also
172 defined “high-confidence” PSI levels as those that meet the following criteria:
173 $\text{*max}(\text{min}(\#C1A, \#AC2), \#C1C2) \geq 5$ AND $\text{min}(\#C1A, \#AC2) + \#C1C2 \geq 10$
174 $\text{*}|\log_2(\#C1A / \#AC2)| \leq 1$ OR $\text{max}(\#C1A, \#AC2) < \#C1C2$

175 For cross-species analyses, we included exons with single-copy orthologues
176 in all species, PSI values in all samples, and confidently alternative spliced in
177 at least one of the samples.

178

179

180

181

Findings

182 ***Data summary***

183 We generated a total of ~909.6 Gb high-quality RNA-seq data, of which ~676.6
184 Gb (~74.6%) reads could reliably aligned to their respective reference genomes
185 (**Supplementary Fig. S3 and Table S1**). We found that on average 61.2%
186 annotated protein coding genes in each genome had FPKM expression values
187 greater than 0.5 (**Supplementary Fig. S4 and Table S3**).

188 ***Concordance in the tree topology based on nucleotide sequence*** 189 ***alignments and gene expression data***

190 Nucleotide alignments-based phylogenetic relationships of these high-altitude
191 vertebrates and their low-altitude relatives matched the established
192 morphological species groupings and the known history of population formation
193 (**Fig. 1b**). The gene expression-based tree based 4,746 transcribed single-copy
194 orthologous genes (66.61% of 7125) for each tissue showed a highly consistent

195 topology to the nucleotide sequence alignment-based phylogeny (**Fig. 2,**
196 **Supplementary Methods**) [9]: mammals were mainly divided into omnivore
197 (pig) and ruminant (goat, sheep and yak/cattle); within the ruminant cluster, the
198 two caprinae (goat and sheep) were closer to each other than the boviniae
199 (yak/cattle). This observation lends supports the idea that gene expression
200 changes evolve together with genetic variation over evolutionary time, resulting
201 in lower expression divergence between more closely species [19].

202 ***Distinctly transcriptomic characteristics between gene expression and***
203 ***alternative splicing***

204 Through comparison of expression levels of 4,746 transcribed single-copy
205 orthologous genes (**Supplementary Fig. S2**) and alternative splicing patterns
206 (reflected by PSI values) of 2,783 orthologous exons shared by the five
207 vertebrates genomes, we observed a tissue-dominated clustering pattern of
208 gene expression, but a species-dominated clustering pattern of alternative
209 splicing [20, 21].

210 For gene expression, there were critical biological differences among tissues
211 (Pearson's $r = 0.67$ and weighted average proportion variance = 0.36), followed
212 by species (Pearson's $r = 0.75$, weighted average proportion variance = 0.22)
213 and local adaptation (Pearson's $r = 0.95$ and weighted average proportion
214 variance = 0.019) (**Fig. 3a** and **Supplementary Fig. S5**). By contrast, for
215 alternative splicing, the differences among species (Pearson's $r = 0.64$ and
216 weighted average proportion variance = 0.30) were higher than among tissues
217 (Pearson's $r = 0.78$ and weighted average proportion variance = 0.075),
218 followed by between high- and low-altitude animals (Pearson's $r = 0.84$ and
219 weighted average proportion variance = 0.021) (**Fig. 3b** and **Supplementary**

220 **Figure S6).**

221 Unsupervised clustering (**Figs. 4a and 4c**) and principal components
222 analysis (PCA) (**Figs. 4b and 4d** and **Supplementary Figs. S7 and S8**) both
223 recapitulated the distinctly transcriptomic characteristics between gene
224 expression and alternative splicing. Tissue-dominated clustering of gene
225 expression indicated that in general tissues possess conserved gene
226 expression signatures and suggested that conserved gene expression
227 differences underlie tissue identity in mammals. On the other hand, greater
228 prominence of species-dominated clustering of alternative splicing suggested
229 that exon splicing is more often affected by species-specific changes in *cis*-
230 regulatory elements and/or *trans*-acting factors than gene expression [20, 21].

231 Notably, tissue-dominated clustering patterns of gene expression further
232 revealed that the cluster of striated muscle (heart and skeletal muscle) and the
233 cluster of vessel-rich tissues (lung and spleen) were closer to each other than
234 the cluster of metabolic tissues (kidney and liver), followed by the distinct
235 clusters of bird (chicken) and mammals according to the evolutionary distance
236 (**Figs. 4a and 4b**). Notably, tissues of birds (chickens) formed a distinct cluster,
237 rather than with their mammalian counterparts, which indicates that divergence
238 in gene expression among these species started to surpass that between
239 different tissues around when birds diverged from mammals (approximately
240 300 million years ago) (**Figs. 4a and 4b**).

241 ***Gene expression plasticity to a high-altitude environment***

242 To exclude the impact of prominence of tissues-dominated clustering of gene
243 expression, so as to comprehensively present transcriptomic differences

244 involved in high-altitude response based on whole annotated genes of their
245 respective genome assembly instead of the single-copy orthologous, we
246 measured the pairwise difference of gene expression between the high-altitude
247 populations and their low-altitude relatives within each tissue for each
248 vertebrate.

249 We identified ~1,423 DEGs between 30 low- versus high-altitude pairs (177
250 DEGs in muscle of chickens to 3,853 DEGs in kidney of sheep) (**Table 1**).
251 Notably, among five pairs of vertebrate, the highly-diverged yak and cattle [1]
252 exhibited the highest number of DEG (~2,005) across six tissues. Among six
253 tissues, the highly aerobic kidney [22] exhibited the highest number of DEGs
254 (~2,097) across five pairs of vertebrates.

255 Expectedly, the more closely related vertebrates (**Fig. 1**) shared more DE
256 genes (**Supplementary Figs. S9–10** and **Additional File 3**). Compared with
257 shared DE genes among mammals, especially between the two closely related
258 members of Caprinae (goat and sheep), the birds (chickens) exhibited
259 significantly fewer shared DE genes with mammals (Wilcoxon rank sum test,
260 $P < 0.0021$) (**Supplementary Fig. S11**). We also identified significantly enriched
261 functional gene categories of DE genes (Chi-square test or Fisher's exact test,
262 $P < 1.03 \times 10^{-4}$), which were shared among multiple pairwise comparisons
263 (**Supplementary Figs. S12–13** and **Additional File 4**), that were potentially
264 related to the dramatic phenotypic changes shaped by high-altitude adaptation,
265 such as response to hypoxia (typically, 'oxidation reduction', 'heme binding',
266 'oxygen binding' , 'oxygen transport' and 'oxygen transporter activity'),
267 cardiovascular system ('angiogenesis' and 'positive regulation of
268 angiogenesis'), the efficiency of biomass production in the resource-poor

269 highland ('metabolic pathways', 'cholesterol biosynthetic process' and 'steroid
270 metabolic process') as well as immune response ('responses of immune and
271 defense') (**Additional file 2**).

272 **Conclusions**

273 High-altitude adaptive evolution of transcription, and the convergence and
274 divergence of transcriptional alteration across species in response to high-
275 altitude environments, is an important topic of broad interest to the general
276 biology community. Here we provide a comprehensive comparative
277 transcriptome landscape of expression and alternative splicing variation
278 between low- and high-altitude populations across multiple species for distinct
279 tissues. Our data serves a valuable resource for further study on gene
280 regulatory changes to adaptive evolution of complex phenotypes.

281 **Availability of supporting data**

282 The RNA-seq data for 180 samples was deposited in the NCBI Gene
283 Expression Omnibus (GEO) under accession numbers GSE93855, GSE77020
284 (Note: GSM1617847-GSM1617849 and GSM2042608-GSM2042610 are
285 duplicates and represent the same samples) and GSE66242 (Note: 9 goat
286 samples derived from individuals sampled at 2000m altitude were not included
287 in this study). The re-sequencing data for 7 individuals was deposited in the
288 NCBI-sequence read archive (SRA) under accession number SRP096151.
289 Supporting data is also available via the *Gigascience* database, GigaDB
290 (GigaDB, RRID:SCR_004002) [23]. Supplementary figures and tables are
291 provided as Additional Files 1-4.

292

293 **Ethics statement**

294 All studies involving animals were conducted according to Regulations for the
295 Administration of Affairs Concerning Experimental Animals (Ministry of Science
296 and Technology, China, revised in June 2004). All experimental procedures and
297 sample collection methods in this study were approved by the Institutional
298 Animal Care and Use Committee of the College of Animal Science and
299 Technology of Sichuan Agricultural University, Sichuan, China, under permit No.
300 DKY-B20121406. Animals were allowed free access to food and water under
301 normal conditions, and were humanely sacrificed as necessary, to ameliorate
302 suffering.

303 **Consent for publication**

304 Not applicable.

305 **Competing interests**

306 The authors declare that they have no competing interests.

307 **Funding**

308 This work was supported by grants from the National High Technology
309 Research and Development Program of China (863 Program) (2013AA102502),
310 the National Natural Science Foundation of China (31402046, 31522055,
311 31601918, 31530073 , 31472081 and 31772576), the Science & Technology
312 Support Program of Sichuan (2016NYZ0042), the Youth Science Fund of
313 Sichuan (2017JQ0011), the China Postdoctoral Science Foundation
314 (2015M572486), China Agriculture Research System (CARS-36), the Program

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

315 for Innovative Research Team of Sichuan Province (2015TD0012), the
316 Program for Pig Industry Technology System Innovation Team of Sichuan
317 Province (SCCXTD-005), the Project of Sichuan Education Department
318 (15ZA0008, 15ZA0003, 16ZA0025 and 16ZB0037), the National Program for
319 Support of Top-notch Young Professionals and the Young Scholars of the
320 Yangtze River.

321 **Authors' contributions**

322 MZ.L., QZ.T., YR.G. and XW.L. designed and supervised the project. JQ.G.,
323 TD.C., SL.H., Y.L., XM.Y., X.T., ZJ.Z., XH.C., DY.L., XL.L. and XB.L. collected
324 the data, L.J., R.L., J.L., KR.L., SL.T., GS.W., JD.M., X.W., MM.M. and AA.J.
325 generated the data. QZ.T. and MZ.L. performed the bioinformatics analyses.
326 QZ.T. and MZ.L. wrote the manuscript. XM.Z. and VN.G. revised the manuscript.

327

328

329

330 **References**

331

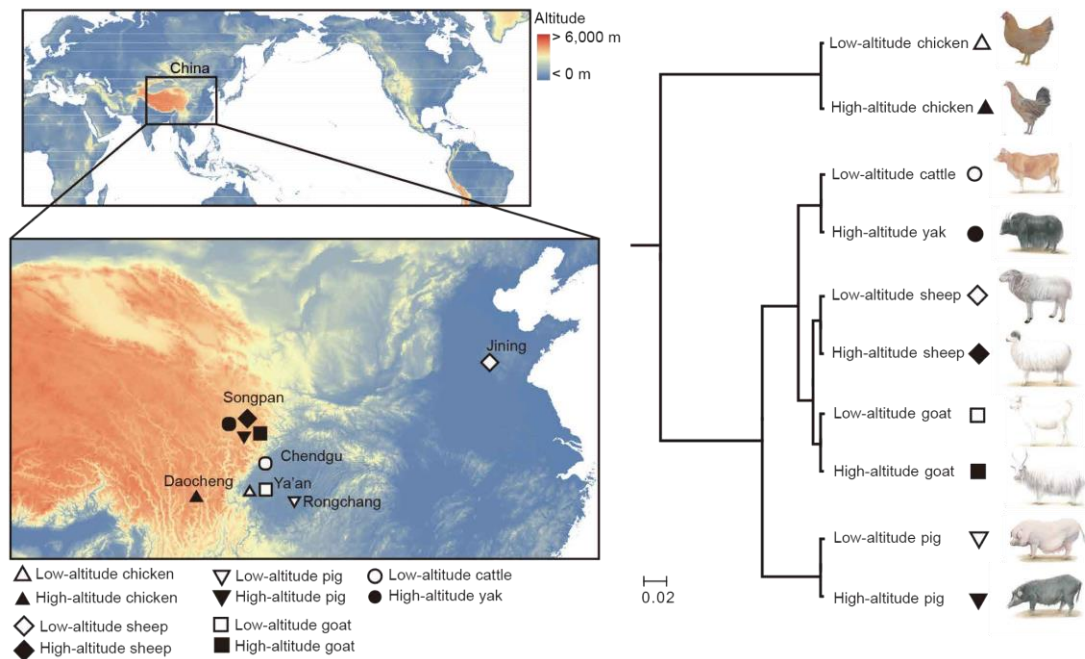
- 332 1. Qiu Q, Zhang G, Ma T, Qian W, Wang J, Ye Z, et al. The yak genome and
333 adaptation to life at high altitude. Nat Genet. 2012;44(8):946-9.
334 doi:10.1038/ng.2343.
- 335 2. Li M, Tian S, Jin L, Zhou G, Li Y, Zhang Y, et al. Genomic analyses identify distinct
336 patterns of selection in domesticated pigs and Tibetan wild boars. Nat Genet.
337 2013;45(12):1431-8. doi:10.1038/ng.2811.
- 338 3. Li M, Chen L, Tian S, Lin Y, Tang Q, Zhou X, et al. Comprehensive variation
339 discovery and recovery of missing sequence in the pig genome using multiple de

1 340 novo assemblies. *Genome Res.* 2017 May;27(5):865-874.
 2 341 doi:10.1101/gr.207456.116.
 3
 4 342 4. International Chicken Genome Sequencing Consortium. Sequence and
 5 comparative analysis of the chicken genome provide unique perspectives on
 6 343 vertebrate evolution. *Nature.* 2004;432(7018):695-716. doi:10.1038/nature03154.
 7 344
 8 345 5. Groenen MA, Archibald AL, Uenishi H, Tuggle CK, Takeuchi Y, Rothschild MF, et
 9 al. Analyses of pig genomes provide insight into porcine demography and evolution.
 10 346 *Nature.* 2012;491(7424):393-8. doi:10.1038/nature11622.
 11 347
 12 348 6. Bovine Genome Sequencing and Analysis Consortium, Elsik CG, Tellam RL,
 13 349 Worley KC, Gibbs RA, et al. The genome sequence of taurine cattle: a window to
 14 350 ruminant biology and evolution. *Science.* 2009;324(5926):522-8.
 15 351 doi:10.1126/science.1169588.
 16 352 7. Dong Y, Xie M, Jiang Y, Xiao N, Du X, Zhang W, et al. Sequencing and automated
 17 353 whole-genome optical mapping of the genome of a domestic goat (*Capra hircus*).
 18 354 *Nat Biotechnol.* 2013;31(2):135-41. doi:10.1038/nbt.2478.
 19 355 8. Jiang Y, Xie M, Chen W, Talbot R, Maddox JF, Faraut T, et al. The sheep genome
 20 356 illuminates biology of the rumen and lipid metabolism. *Science.* 2014;344
 21 357 (6188):1168-73. doi:10.1126/science.1252806.
 22 358 9. Vilella AJ, Severin J, Ureta-Vidal A, Heng L, Durbin R and Birney E.
 23 359 EnsemblCompara GeneTrees: Complete, duplication-aware phylogenetic trees in
 24 360 vertebrates. *Genome Res.* 2009;19(2):327-35. doi:10.1101/gr.073585.107.
 25 361 10. Li H and Durbin R. Fast and accurate long-read alignment with Burrows-Wheeler
 26 362 transform. *Bioinformatics.* 2010;26(5):589-95. doi:10.1093/bioinformatics/btp698.
 27 363 11. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernysky A, et al.
 28 364 The Genome Analysis Toolkit: a MapReduce framework for analyzing next-
 29 365 generation DNA sequencing data. *Genome Res.* 2010;20(9):1297-303.
 30 366 doi:10.1101/gr.107524.110.
 31 367 12. Li H, Coghlan A, Ruan J, Coin LJ, Heriche JK, Osmotherly L, et al. TreeFam: a

1 368 curated database of phylogenetic trees of animal gene families. *Nucleic Acids*
2 369 *Res.* 2006;34(Database issue):D572-80. doi:10.1093/nar/gkj118.
3
4 370 13. Trapnell C, Pachter L and Salzberg SL. TopHat: discovering splice junctions with
5 RNA-Seq. *Bioinformatics.* 2009;25(9):1105-11. doi:10.1093/bioinformatics/btp120.
6 371
7 372 14. Trapnell C, Williams BA, Pertea G, Mortazavi A, Kwan G, van Baren MJ, et al.
8 Transcript assembly and quantification by RNA-Seq reveals unannotated
9 373 transcripts and isoform switching during cell differentiation. *Nat Biotechnol.*
10 374 2010;28(5):511-5. doi:10.1038/nbt.1621.
11 375
12 376 15. The pvca R package.
13 <https://bioconductor.org/packages/release/bioc/html/pvca.html>. Accessed Feb 16
14 377 2017.
15 378
16 379 16. Paradis E, Claude J and Strimmer K. APE: Analyses of phylogenetics and
17 380 evolution in R language. *Bioinformatics.* 2004;20(2):289-90.
18 381
19 382 17. Dennis G, Jr., Sherman BT, Hosack DA, Yang J, Gao W, Lane HC, et al. DAVID:
20 383 Database for annotation, visualization, and integrated discovery. *Genome Biol.*
21 2003;4(5):P3.
22 384
23 385 18. Kent WJ, Sugnet CW, Furey TS, Roskin KM, Pringle TH, Zahler AM, et al. The
24 386 human genome browser at UCSC. *Genome Res.* 2002;12(6):996-1006.
25 387 doi:10.1101/gr.229102.
26 388
27 389 19. Brawand D, Soumillon M, Necsulea A, Julien P, Csardi G, Harrigan P, et al. The
28 390 evolution of gene expression levels in mammalian organs. *Nature.* 2011;478
29 391 (7369):343-8. doi:10.1038/nature10532.
30 392
31 393 20. Merkin J, Russell C, Chen P and Burge CB. Evolutionary dynamics of gene and
32 394 isoform regulation in Mammalian tissues. *Science.* 2012;338(6114):1593-9.
33 395 doi:10.1126/science.1228186.
34
35 396 21. Barbosa-Morais NL, Irimia M, Pan Q, Xiong HY, Gueroussov S, Lee LJ, et al. The
36 397 evolutionary landscape of alternative splicing in vertebrate species. *Science.*
37 2012;338(6114):1587-93. doi:10.1126/science.1230612.
38 398
39 399
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

- 396 22. Mele M, Ferreira PG, Reverter F, DeLuca DS, Monlong J, Sammeth M, et al.
 397 Human genomics. The human transcriptome across tissues and individuals.
 398 Science. 2015;348(6235):660-5. doi:10.1126/science.aaa0355.
 399 23. Tang Q, Gu Y, Zhou X, Jin L, Guan J, Liu R, et al. Supporting data for "Comparative
 400 transcriptomics of five high-altitude vertebrates and their low-altitude relatives".
 401 *GigaScience* Database. 2017. <http://dx.doi.org/10.5524/100355>

408 **Figures 1-4**



410
411
412 **Figure 1. Sampling locations and nucleotide alignment-based tree.**

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

413 (a) Geographic locations of the studied animals.

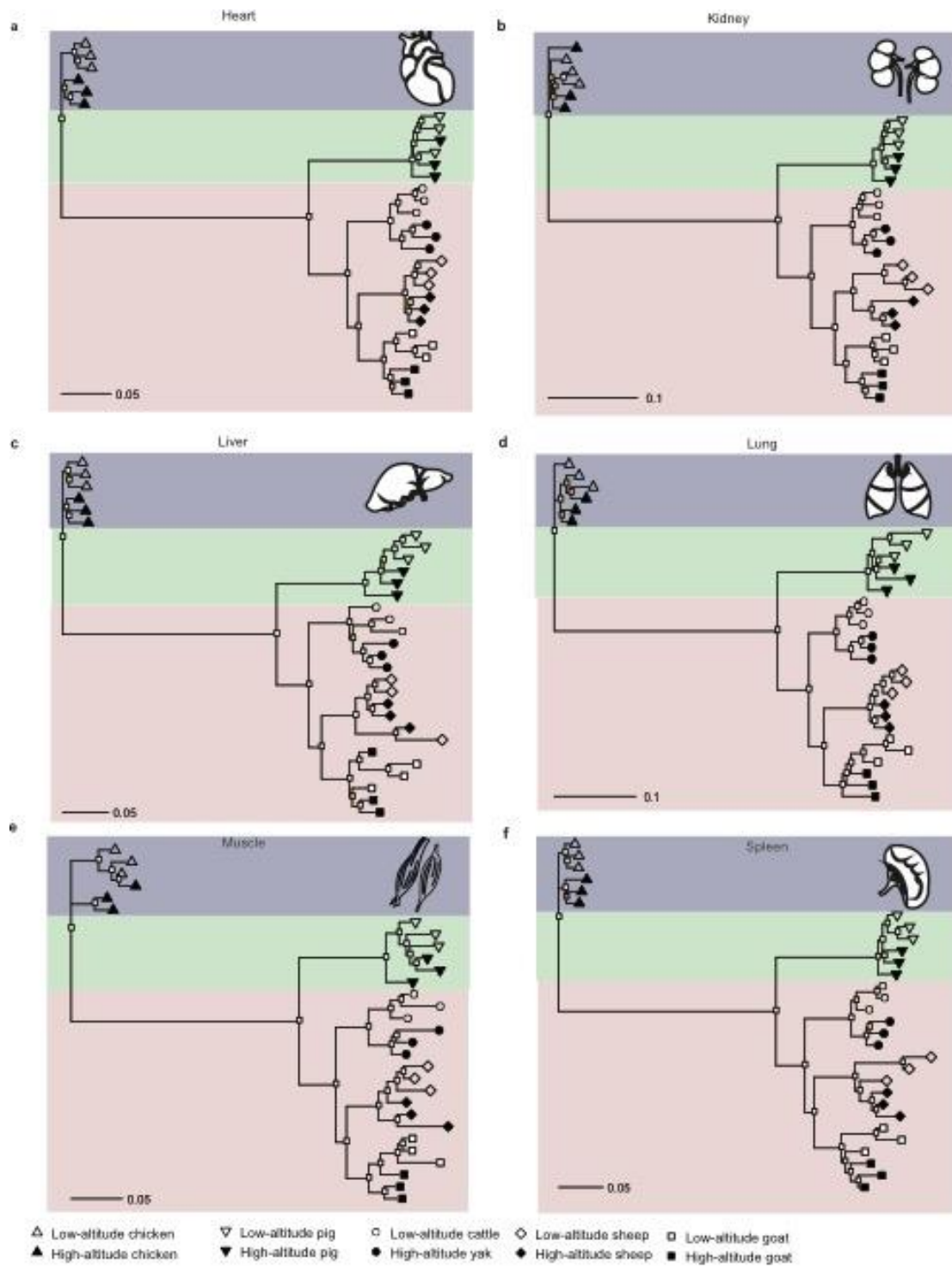
414 (b) A neighbour-joining tree constructed based on concatenated coding sequences of
415 single-copy orthologues substituted by SNVs and InDels detected in each animal.

416 We downloaded and extracted the unassembled reads from short-insert (500 bp) libraries
417 of a single yak [1], a Tibetan pig [2] and a Rongchang pig [3] that were used for *de novo*
418 assemblies to roughly 10 × depth coverage. We also randomly selected an individual of
419 the cattle, low- and high-altitude chicken, goat and sheep and sequenced the whole
420 genomes at ~10 × depth coverage.

421

422

423



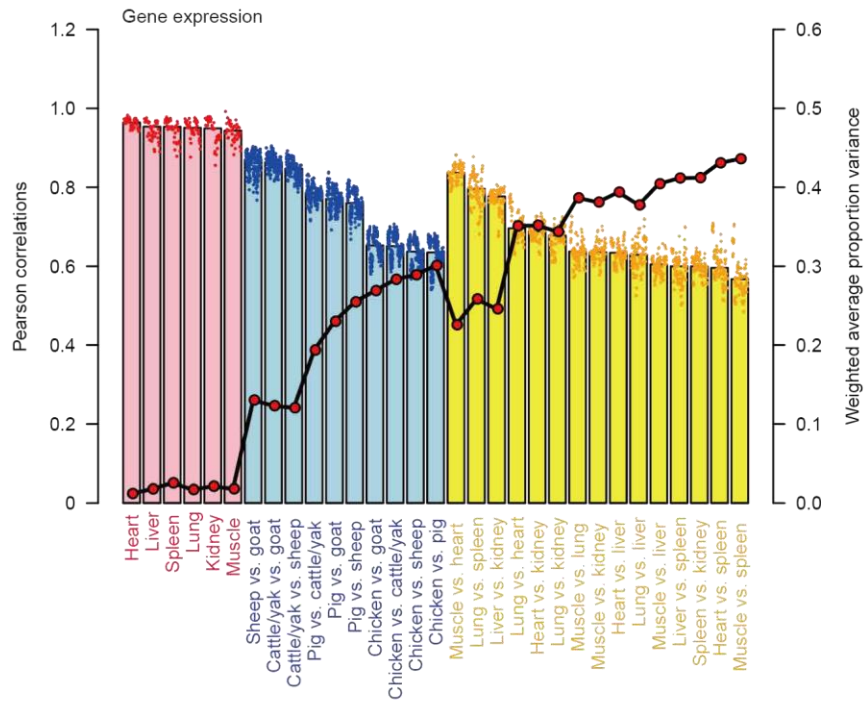
424

425 **Figure 2. Gene expression phylogenies for six tissues across five vertebrates.**

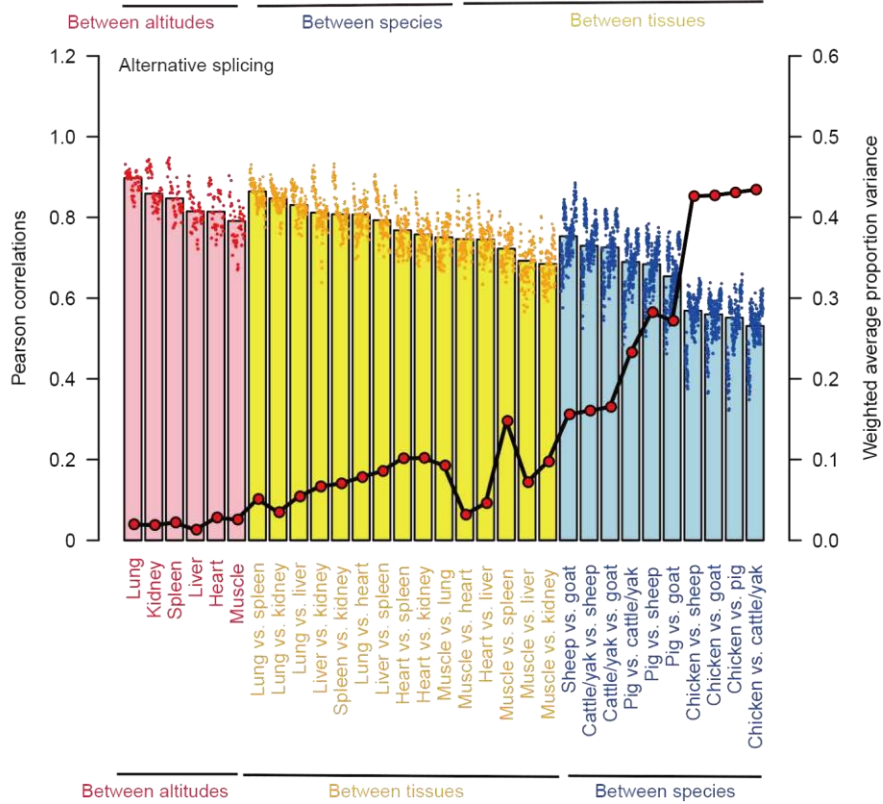
426 Neighbour-joining expression tree constructed based on (1-Spearman correlation)
 427 distances in six tissues. We performed 100 bootstraps by randomly sampling the single
 428 copy orthologues with replacement. Bootstrap values (fractions of replicate trees that have
 429 the branching pattern as in the shown tree constructed using all the transcribed single copy
 430 orthologues) are indicated by different colors: red color of the node indicates support from

431 less than 50% bootstraps, while orange, yellow and white colors indicate support between
 432 50% and 70%, between 70% and 90% and more than 90%, respectively.

a

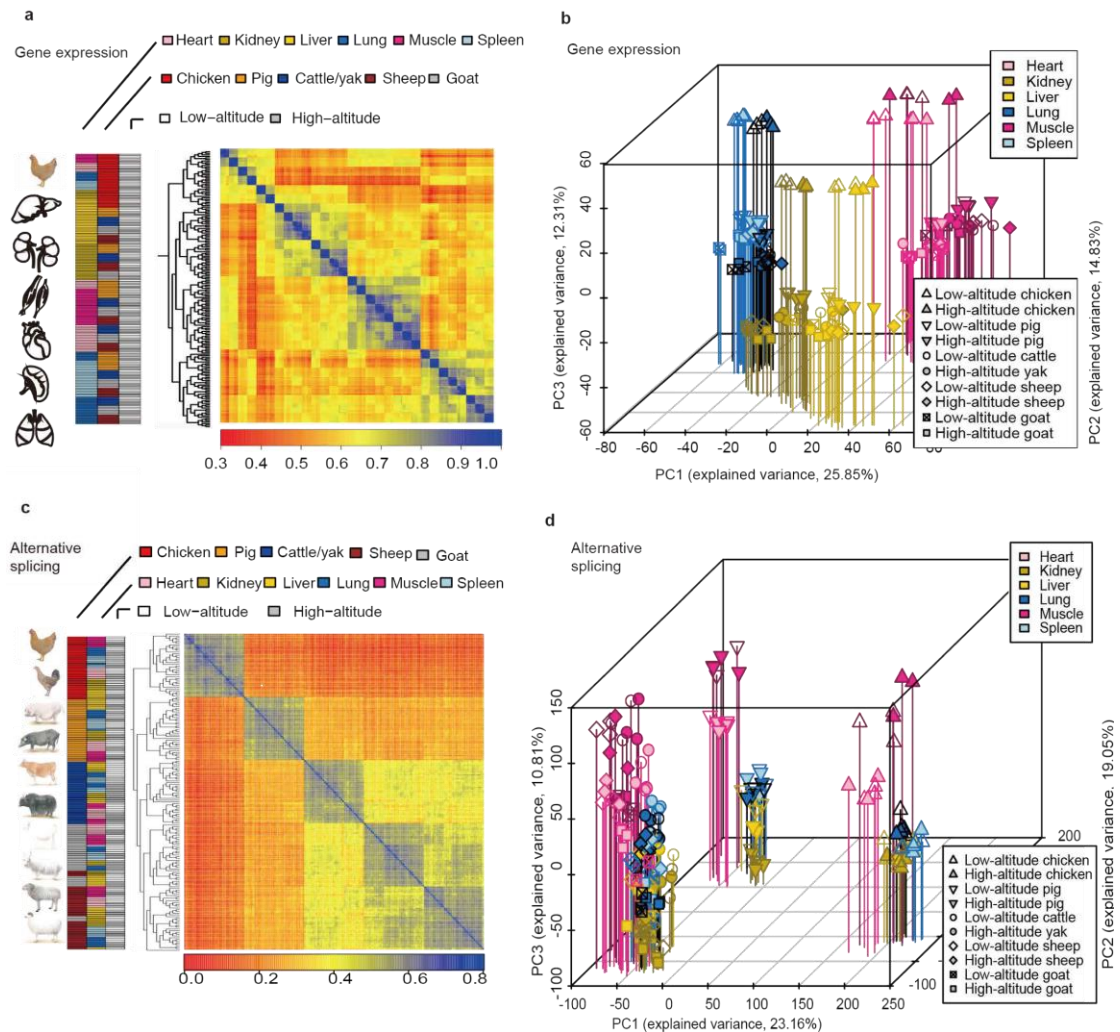


b



433
 434 **Figure 3. Comparison of variations between altitudes, species and tissues revealed**
 435 **by (a) gene expression and (b) alternative splicing pattern.**

436 Scatter-point and bar plots represent the pairwise Pearson's correlation between samples.
 437 Weighted average proportion variance of the alternative splicing (reflected by PSI values)
 438 were determined using the Principal Variance Component Analysis (PVCA) approach and
 439 depicted as red dots connected by black lines.



440

441 **Figure 4. Global pattern of gene expression and alternative splicing pattern.**

442 Hierarchical clustering of samples using (a) the gene expression and (c) the alternative
 443 splicing (reflected by PSI values). Average linkage hierarchical clustering was used with
 444 distance between samples measured by the Pearson's correlation between the vectors of
 445 expression values. Factorial map of the principal-component analysis (PCA) of (b) gene

1 446 expression levels and **(d)** the alternative splicing. The proportion of the variance
2
3
4 447 explained by the principal components is indicated in parentheses. The vertical leading
5
6 448 lines with different colors from the plotted points dropping to the xy-plane show the
7
8
9 449 separation of points based on the first and second principal components.
10
11 450
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

Table 1. Number of DEGs between five high-altitude vertebrates and their low-altitude relatives for each tissue

Species	Heart	Kidney	Liver	Lung	Muscle	Spleen	Mean
Chicken	1283 (8.28%)	748 (4.83%)	613 (3.96%)	1072 (6.92%)	177 (1.14%)	984 (6.35%)	812 (5.25%)
Pig	206 (0.95%)	532 (2.46%)	1199 (5.55%)	426 (1.97%)	385 (1.78%)	994 (4.60%)	623 (2.89%)
Cattle/yak	1602 (8.02%)	1797 (8.99%)	869 (4.35%)	3092 (15.47%)	2403 (12.03%)	2268 (11.35%)	2005 (10.04%)
Sheep	1332 (6.37%)	3853 (18.43%)	259 (1.24%)	1829 (8.75%)	1079 (5.16%)	2356 (11.27%)	1784 (8.54%)
Goat	2215 (10.01%)	3557 (16.07%)	655 (2.96%)	1330 (6.01%)	2305 (10.42%)	1269 (5.73%)	1888 (8.53%)
Mean	1327 (6.73%)	2097 (10.16%)	719 (3.61%)	1549 (7.82%)	1269 (6.11%)	1574 (7.86%)	

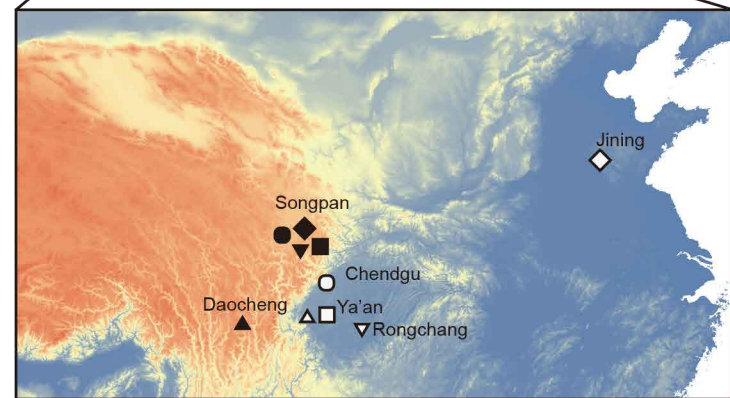
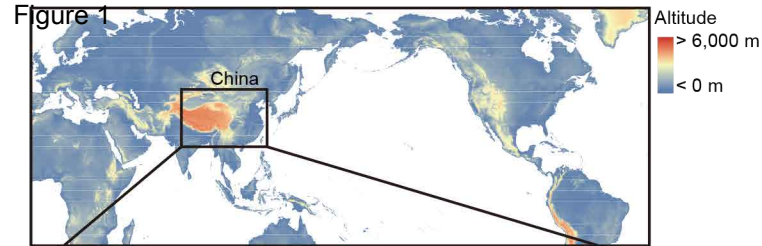
Percentage of the DGEs compared with the total number of annotated protein coding genes in their respective reference genomes are listed in parenthesis. There are 15495, 21594, 19981, 22131, 20908 annotated protein coding genes in reference genomes of Chicken (Galgal4) [4], pig (Suscrofa 10.2) [5], cattle (UMD3.1) [6], goat (CHIR_1.0) [7] and sheep (Oar_v3.1) [8], respectively.

Table 1. Number of DEGs between five high-altitude vertebrates and their low-altitude relatives for each tissue

Species	Heart	Kidney	Liver	Lung	Muscle	Spleen	Mean
Chicken	1283 (8.28%)	748 (4.83%)	613 (3.96%)	1072 (6.92%)	177 (1.14%)	984 (6.35%)	812 (5.25%)
Pig	206 (0.95%)	532 (2.46%)	1199 (5.55%)	426 (1.97%)	385 (1.78%)	994 (4.60%)	623 (2.89%)
Cattle/yak	1602 (8.02%)	1797 (8.99%)	869 (4.35%)	3092 (15.47%)	2403 (12.03%)	2268 (11.35%)	2005 (10.04%)
Sheep	1332 (6.37%)	3853 (18.43%)	259 (1.24%)	1829 (8.75%)	1079 (5.16%)	2356 (11.27%)	1784 (8.54%)
Goat	2215 (10.01%)	3557 (16.07%)	655 (2.96%)	1330 (6.01%)	2305 (10.42%)	1269 (5.73%)	1888 (8.53%)
Mean	1327 (6.73%)	2097 (10.16%)	719 (3.61%)	1549 (7.82%)	1269 (6.11%)	1574 (7.86%)	

Percentage of the DGEs compared with the total number of annotated protein coding genes in their respective reference genomes are listed in parenthesis. There are 15495, 21594, 19981, 22131, 20908 annotated protein coding genes in reference genomes of Chicken (Galgal4) [4], pig (Suscrofa 10.2) [5], cattle (UMD3.1) [6], goat (CHIR_1.0) [7] and sheep (Oar_v3.1) [8], respectively.

Figure 1



- △ Low-altitude chicken
- ▲ High-altitude chicken
- ◇ Low-altitude sheep
- ◆ High-altitude sheep
- ▽ Low-altitude pig
- ▼ High-altitude pig
- Low-altitude cattle
- High-altitude yak
- Low-altitude goat
- High-altitude goat

[Click here to download Figure Figure1.pdf](#)

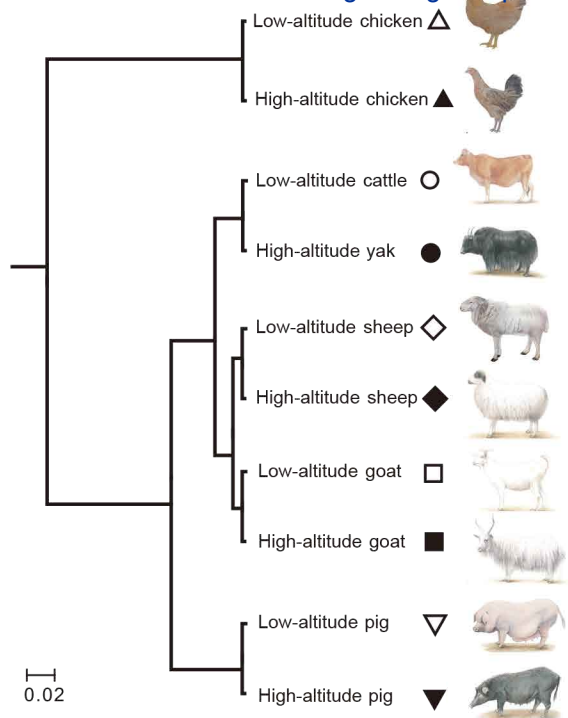
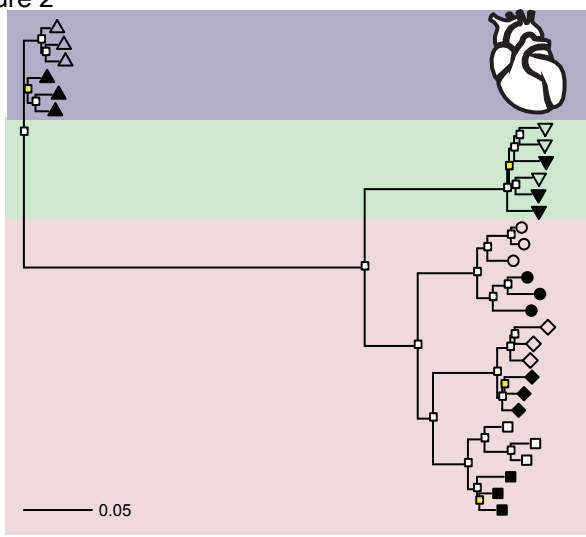
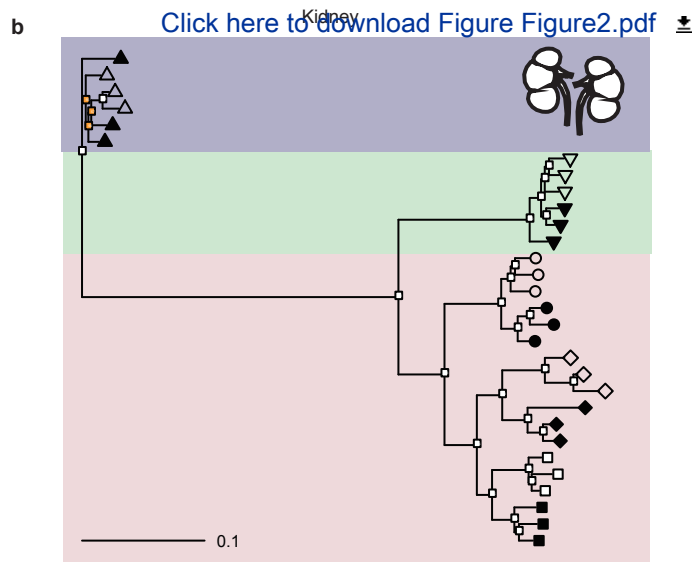


Figure 2

Heart

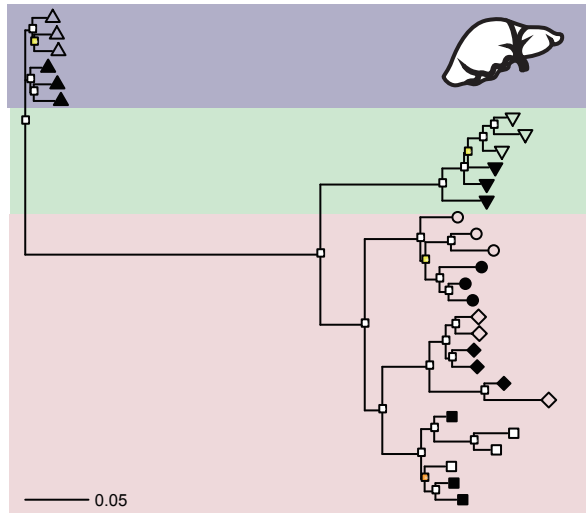


Kidney

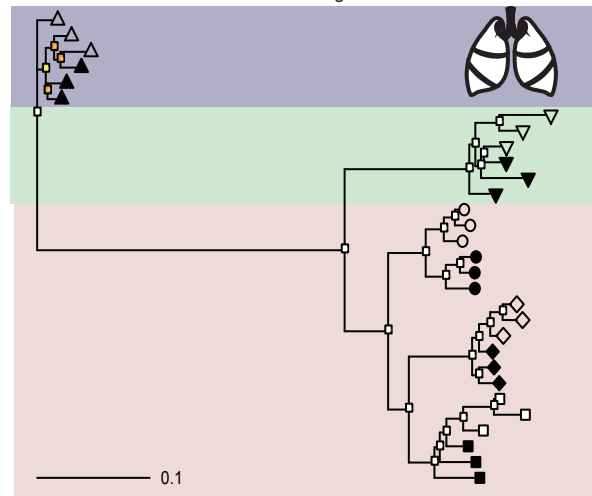


[Click here to download Figure2.pdf](#)

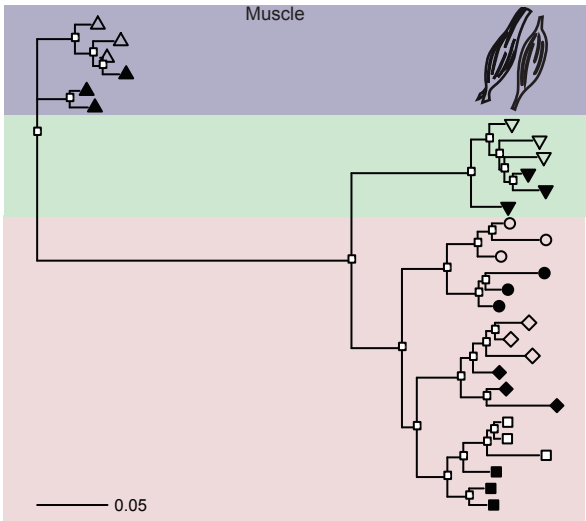
c Liver



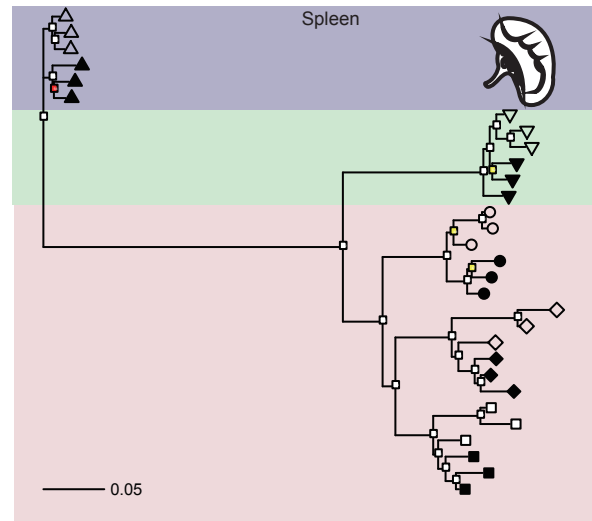
d Lung



e Muscle



f Spleen



△ Low-altitude chicken
▲ High-altitude chicken

▽ Low-altitude pig
▼ High-altitude pig

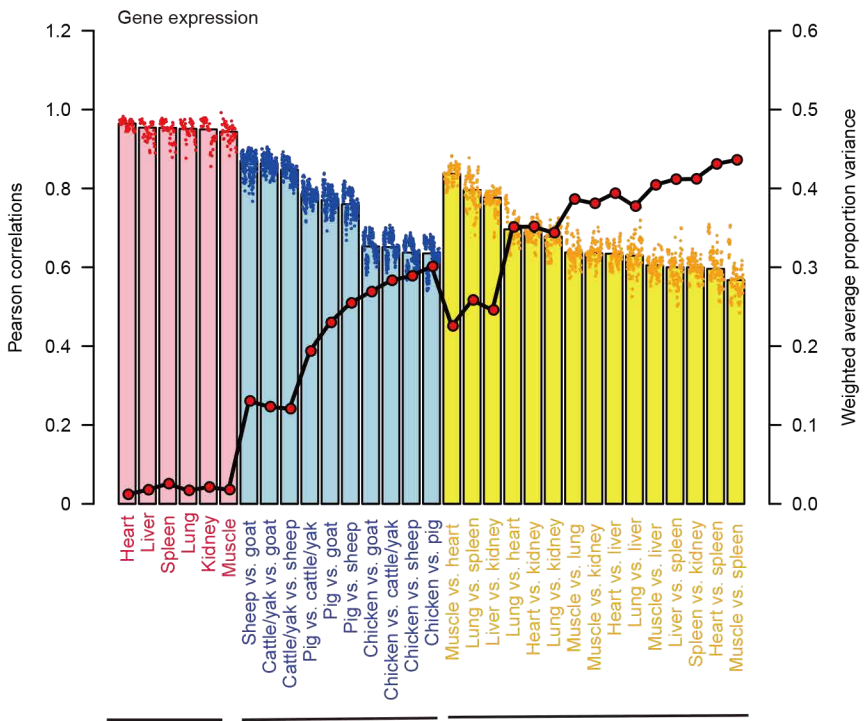
○ Low-altitude cattle
● High-altitude yak

◇ Low-altitude sheep
◆ High-altitude sheep

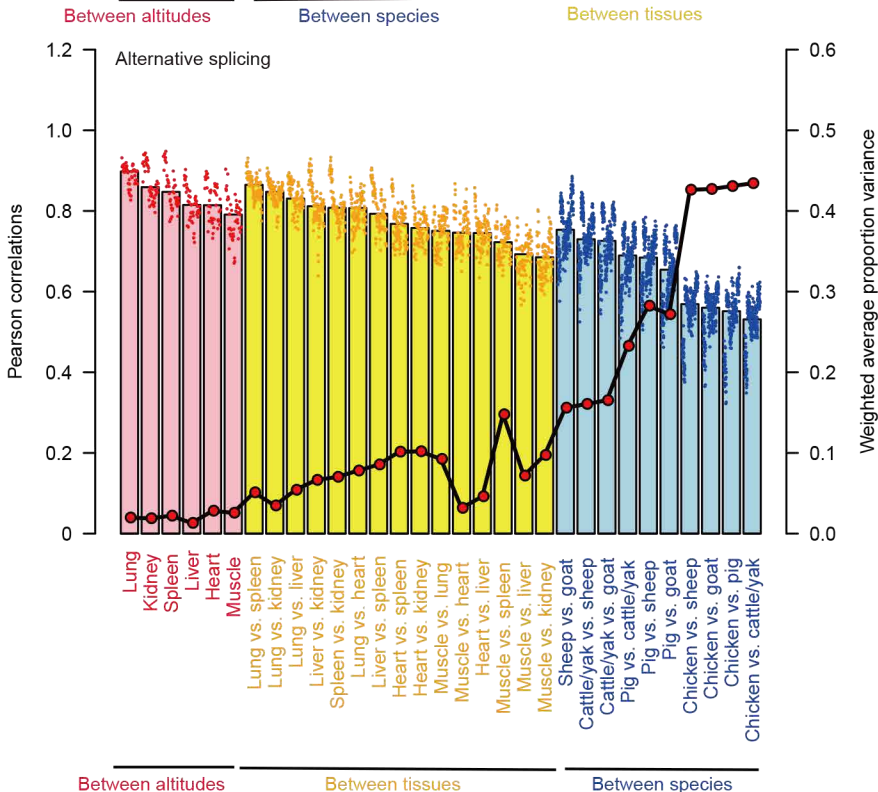
□ Low-altitude goat
■ High-altitude goat

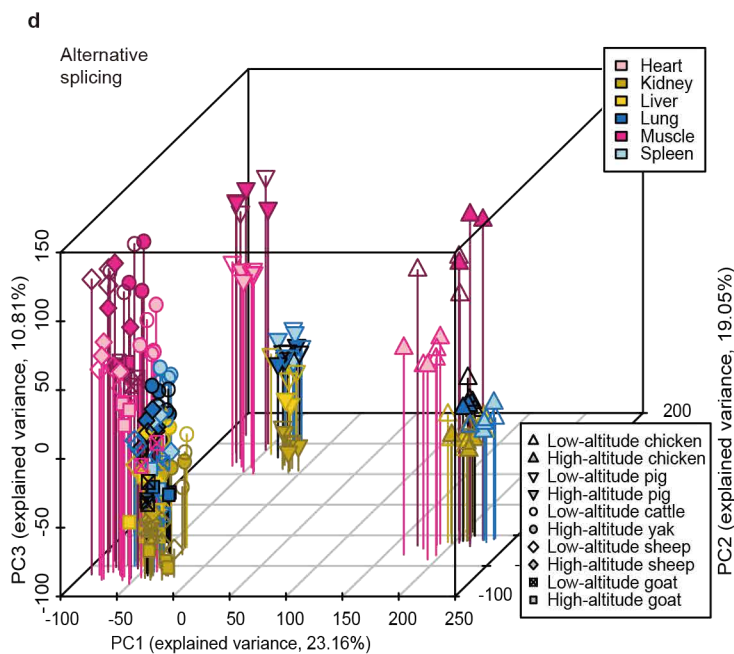
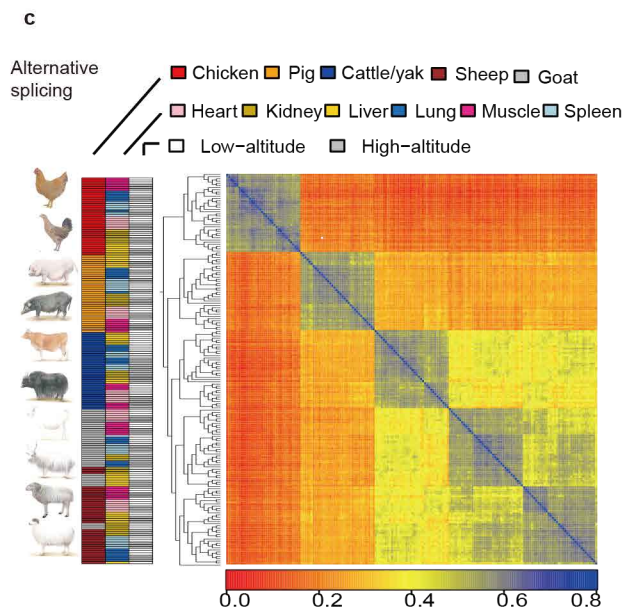
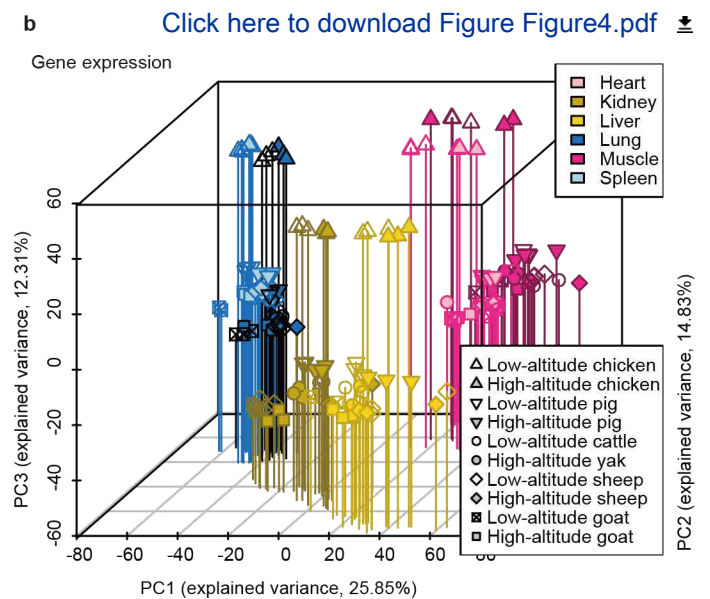
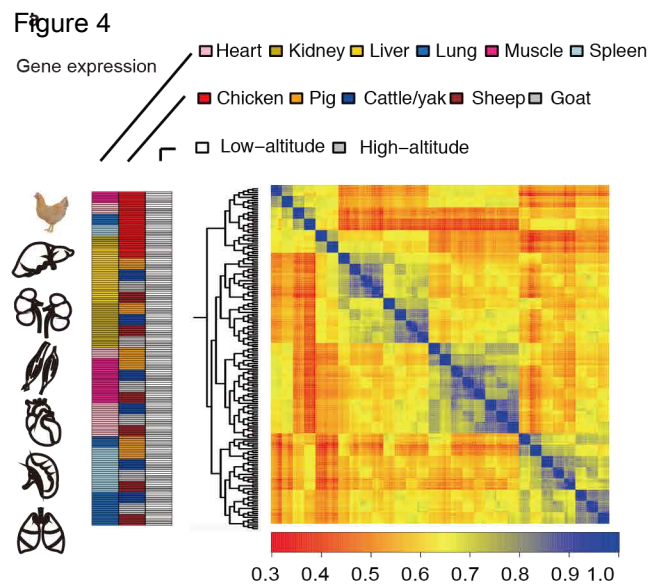
Figure 3

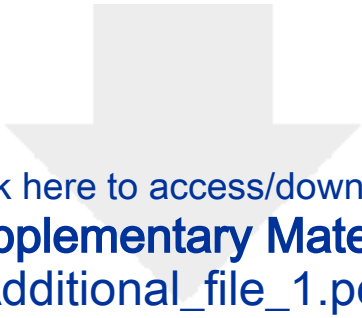
[Click here to download Figure Figure3.pdf](#)




b

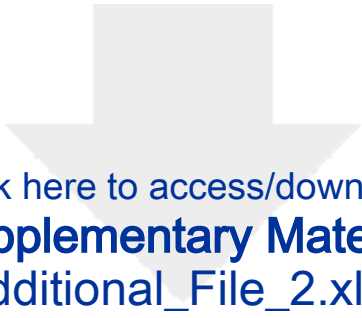







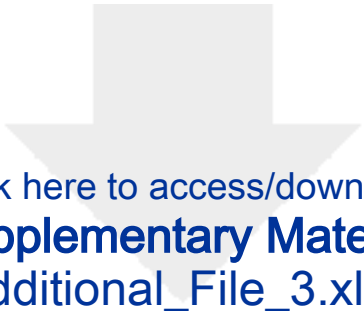
Click here to access/download
Supplementary Material
Additional_file_1.pdf






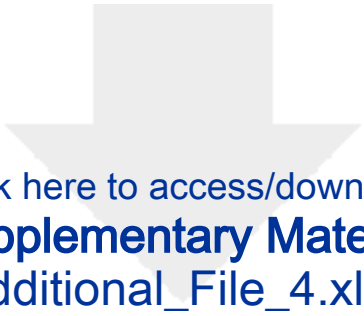
Click here to access/download
Supplementary Material
Additional_File_2.xlsx






Click here to access/download
Supplementary Material
Additional_File_3.xlsx





Click here to access/download
Supplementary Material
Additional_File_4.xlsx





Click here to access/download
Supplementary Material
Response_sup.pdf



GigaScience

em@editorialmanager.com

Dear Dr. Hans Zauner,

We are delighted to be informed of the positive responses from you. We have carefully revised the manuscript following your and review 2's suggestions, and uploaded it online for formal acceptance and publication.

We have released all of datasets related to our manuscript, including a total of 180 RNA-seq data deposited in the NCBI Gene Expression Omnibus (GEO) under accession numbers GSE93855, GSE77020 and GSE66242, as well as the 7 whole-genome sequencing data deposited in the NCBI-sequence read archive (SRA) under accession number SRP096151.

Moreover, to promote the clarity of depiction for the principal component analysis (PCA) results, we newly added the two-dimensional PCA figures as the **Supplementary Figs. S14-S15**.

We sincerely appreciate your assistance in improving the manuscript. We are glad to be able to contribute to *GigaScience*.

Best regards,

Dr. Mingzhou Li

Sichuan Agricultural University, Chengdu, Sichuan, China

Email: mingzhou.li@sicau.edu.cn

Detailed responses to editor

All comments provided by editor are in gray italics, and our responses are in black. Important revisions in the manuscript are marked in red.

Editor: Dr. Hans Zauner

Comment 1-1:

Before we proceed to acceptance, please address the minor additional comment of reviewer 2 (regarding Fig. 4 - see below).

Formatted: Justified, Line spacing: 1.5 lines

Response 1-1:

Thank you for your reminder. We have carefully revised **Fig. 4** according to reviewer 2's suggestion. (Please see **Response 2-3** for details)

Comment 1-2:

In preparation for publication, please also remove all highlighting/tracking that was added for the purpose of review. Please also be aware that prior to publication, all data has to be openly available (e.g. via NCBI - please lift any embargoes etc.).

Response 1-2:

Thank you for your reminder. We have removed all highlighting/tracking that was added for the purpose of review and prepared the manuscripts for publication.

We have released all of dataset related to our manuscript, inducing a total of 180 RNA-seq data deposited in the NCBI Gene Expression Omnibus (GEO) under accession numbers GSE93855, GSE77020 and GSE66242, as well as

7 whole-genome sequencing data deposited in the NCBI-sequence read archive (SRA) under accession number SRP096151.

GSE93855: <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE93855>

GSE77020: <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE77020>

GSE66242: <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE66242>

SRP096151: <https://trace.ncbi.nlm.nih.gov/Traces/sra/?study=SRP096151>

Comment 1-3:

I note you indicate several corresponding authors - please note that we can consider only one author with this role. Please discuss this with your coauthors and indicate only one corresponding author in your revised manuscript.

Response 1-3:

Thank you for your reminder. We have thoroughly discussed the corresponding authorship issue with all of our co-authors, and only indicated two corresponding authors and three equally contributing first authors in the revised manuscript.