

# GigaScience

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

--Manuscript Draft--

<b>Manuscript Number:</b>	GIGA-D-17-00037R1	
<b>Full Title:</b>	Comparative transcriptomics of five high-altitude vertebrates and their low-altitude relatives	
<b>Article Type:</b>	Data Note	
<b>Funding Information:</b>	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
<b>Abstract:</b>	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 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>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>
<b>Corresponding Author:</b>	Mingzhou Li, Ph.D. Sichuan Agricultural University Chengdu, Sichuan CHINA
<b>Corresponding Author Secondary Information:</b>	
<b>Corresponding Author's Institution:</b>	Sichuan Agricultural University
<b>Corresponding Author's Secondary Institution:</b>	
<b>First Author:</b>	Qianzi Tang
<b>First Author Secondary Information:</b>	
<b>Order of Authors:</b>	<p>Qianzi Tang</p> <p>Yiren Gu</p> <p>Xuming Zhou</p> <p>Long Jin</p> <p>Jiuqiang Guan</p> <p>Rui Liu</p> <p>Jing Li</p> <p>Keren Long</p> <p>Shilin Tian</p> <p>Tiandong Che</p> <p>Silu Hu</p> <p>Yan Liang</p> <p>Xuemei Yang</p> <p>Xuan Tao</p> <p>Zhijun Zhong</p> <p>Guosong Wang</p> <p>Xiaohui Chen</p> <p>Diyan Li</p> <p>Jideng Ma</p>

	Xun Wang
	Miaomiao Mai
	An'an Jiang
	Xiaolin Luo
	Xuebin Lv
	Vadim N. Gladyshev
	Xuewei Li
	Mingzhou Li, Ph.D.
<b>Order of Authors Secondary Information:</b>	
<b>Response to Reviewers:</b>	<p>Reviewer 1:</p> <p>Comment 2-1 The authors report a well-developed project to better understand the gene expression differences in multiple tissues from 5 species (with the cattle-yak comparison counted as one). The data collected is enormous and clearly appears to be sufficient for the analyses proposed, but there are a number of questions regarding both the methods and results presented.</p> <p>Response 2-1 Thank you for your positive comments. We would fully address your concerns and provide our point-to-point responses as follows.</p> <p>Comment 2-2 Of highest importance, the authors present a set of analyses in which the output is a list of genes and a calculated expression level; these lists are then used in a number of ways to calculate expression and enriched function per tissue in several comparison. These lists (and not even the numbers of genes in each list) are not provided, so it is impossible to see these lists or use the lists as a resource for other work. Since one aspect of this publication would be as a resource for others, the authors must provide these lists as well as the calculated expression value for each gene. I realize these lists are extensive, but are a crucial component of the resource, especially for those readers who will not start with the raw data, but also for those who can repeat the analyses and compare their resulting normalized expression data with those that the authors created.</p> <p>Response 2-2 Thank you for your reminder. According to the submission guidelines of GigaScience, we uploaded the complete gene lists with normalized expression values to the GigaScience temporary FTP server.</p> <p>Comment 2-3 Further, the authors describe some biological results on comparisons between high and low altitude, but fail to provide sufficient description of the results. The Supplementary file is incomplete (see below), but also the text on all tissue and species comparisons is only a few sentences. More is needed to justify this reporting. For example, a strength of the work is the multi-species comparison of the same question of adaptation to high altitude. A comparison of the high/low differentially expressed gene lists in the same tissue across species would seem minimal and potentially very interesting- i.e., are the genes and pathways identified similar (more similar than expectation?). This would provide more insight, as well as more evidence the analyses are providing biologically relevant information.</p> <p>Response 2-3 Thank you for your valuable suggestions. Based on your suggestions, we evaluated the amount of shared DE genes between the high- and low-altitude populations in each tissue among five vertebrates (Supplemental Figs. S9–10 and Additional File 3), and found that more closely related vertebrates shared more common DE genes (Supplemental Fig. S11). We also discovered that the enriched functional categories of DE genes substantially overlapped (Supplemental Figs. S12–13 and Additional File 4).</p>

We added Supplemental Figs. S9–13 and Additional Files 3–4 to the manuscript.

As shown in the newly added Supplemental Figs. S9-13 and Additional Files 3-4, expectedly, the more closely related vertebrates (Fig. 1) shared more DE genes (Supplementary Figs. S9–10 and Additional File 3). Compared with shared DE genes among mammals, especially between the two closely related members of Caprinae (goat and sheep), the birds (chickens) exhibited significantly fewer shared DE genes with mammals (Wilcoxon rank sum test,  $P < 0.0021$ ) (Supplementary Fig. S11). We also identified significantly enriched functional gene categories of DE genes (Chi-square test or Fisher's exact test,  $P < 1.03 \times 10^{-4}$ ), which were shared among multiple pairwise comparisons (Supplementary Figs. S12–13 and Additional File 4), that were potentially related to the dramatic phenotypic changes shaped by high-altitude adaptation, such as response to hypoxia (typically, 'oxidation reduction', 'heme binding', 'oxygen binding', 'oxygen transport' and 'oxygen transporter activity'), cardiovascular system ('angiogenesis' and 'positive regulation of angiogenesis'), the efficiency of biomass production in the resource-poor highland ('metabolic pathways', 'cholesterol biosynthetic process' and 'steroid metabolic process') as well as immune response ('responses of immune and defense') (Additional file 2) (the statement has been added to the main text, page 11, line 251-267).

#### Comment 2-4

##### 1. Criterion for expression.

a) On line 40, the authors indicate they are using a FPKM of 0.1. I was unable to find specific details on the sequencing data so that I could determine the number of counts this represents. I could not find the read length nor whether this was SE or PE.

Assuming 100 nt read length and PE for the average of 5 Gb for each tissue reported, a FPKM of 0.1 is 2.5 counts for a 1 kb transcript. This is very low. The authors should justify this low cutoff, which affects all subsequent analyses. I would like to see the median expression level for each tissue, as well.

#### Response 2-4

Thank you for your valuable suggestions. Our data are paired-end reads of 100 nt for three tissues (heart, lung, and muscle), and 125 nt for the other three tissues (kidney, liver, and spleen). Although some previous reports used FPKM  $> 0.1$  as the cutoff for transcribed genes [1-3], based on your suggestions, we used a stricter cut-off of FPKM  $> 0.5$  ( $> 0.5$  FPKM for over 80% of the samples) in the subsequent analyses and updated all of the figures and tables. Our findings did not conflict with those in the initial manuscript, and were further strengthened, typically the 3D PCA result: chickens formed a distinct cluster from the mammals, which indicates that divergence in gene expression among these species started to surpass that between different tissues around when birds diverged from mammals (approximately 300 million years). We revised the corresponding text from "The exceptions to tissue dominance were that chicken heart, lung and liver clustered with chicken skeletal muscle, spleen and kidney, respectively, rather than with their mammalian counterparts, which implied that divergence in gene expression among these species started to surpass those between different tissues at about the time when birds split from mammals (~300 million years)" to "Notably, tissues of birds (chickens) formed a distinct cluster, rather than with their mammalian counterparts, which indicates that divergence in gene expression among these species started to surpass that between different tissues around when birds diverged from mammals (approximately 300 million years ago)." (Main text, page 10, lines 232-236). After adding the FPKM 0.5 cut-off filtering for genes and 5 as the gene number cut-off for enriched terms, some of the specific over-represented terms changed even though the enriched general categories remained unchanged. We have revised the corresponding text from "As expected, respectable significantly enriched functional gene categories by DGEs, which shared in multiple pair-wise comparisons, were potentially related to the dramatic phenotypic changes shaped by high-altitude adaptation, such as response to hypoxia (typically, 'oxidation reduction', 'heme binding', 'oxygen binding', 'response to oxygen levels' and 'response to hypoxia'), cardiovascular system ('blood vessel development', 'blood vessel morphogenesis', 'blood circulation' and 'development of lung and heart'), the efficiency of biomass production in the resource-poor highland (processes of 'steroid biosynthesis' and 'fatty acid metabolism') as well as immune response ('responses of immune and defense')" to "Expectedly, the more closely related vertebrates (Fig. 1) shared more DE genes (Supplementary Figs. S9–10 and Additional File 3). Compared with shared DE genes among mammals, especially between the two closely related members of Caprinae



(goat and sheep), the birds (chickens) exhibited significantly fewer shared DE genes with mammals (Wilcoxon rank sum test,  $P < 0.0021$ ) (Supplementary Fig. S11). We also identified significantly enriched functional gene categories of DE genes (Chi-square test or Fisher's exact test,  $P < 1.03 \times 10^{-4}$ ), which were shared among multiple pairwise comparisons (Supplementary Figs. S12–13 and Additional File 4), that were potentially related to the dramatic phenotypic changes shaped by high-altitude adaptation, such as response to hypoxia (typically, 'oxidation reduction', 'heme binding', 'oxygen binding', 'oxygen transport' and 'oxygen transporter activity'), cardiovascular system ('angiogenesis' and 'positive regulation of angiogenesis'), the efficiency of biomass production in the resource-poor highland ('metabolic pathways', 'cholesterol biosynthetic process' and 'steroid metabolic process') as well as immune response ('responses of immune and defense') (Additional file 2)." (Main text, page 11, lines 251-267). We also revised the corresponding text from "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" to "Of these, ~75% reads could be aligned to their respective reference genomes, and on average ~60% of annotated protein coding genes in each organism showed FPKM expression values greater than 0.5" (Main text, page 2, lines 40-41); from "Log2-transformed values of (FPKM + 1) for genes were used in subsequent analyses" to "Log2-transformed values of (FPKM + 1) for genes with >0.5 FPKM in over 80% of the samples were used in subsequent analyses" (Main text, page 5, lines 113-114); from "We found that on average 69.7% annotated protein coding genes in each genome had FPKM expression values greater than 0.1" to "We found that on average 61.2% annotated protein coding genes in each genome had FPKM expression values greater than 0.5" (Main text, page 8, lines 181-183); from "The gene expression-based tree based 7,125 single-copy orthologous genes for each tissue showed a highly consistent topology to the nucleotide sequence alignment-based phylogeny" to "The gene expression-based tree based 4,746 transcribed single-copy orthologous genes (66.61% of 7125) for each tissue showed a highly consistent topology to the nucleotide sequence alignment-based phylogeny (Fig. 2, Supplementary Methods) [9]" (Main text, page 8, lines 189-192); from "Through comparison of expression levels of 7,125 single-copy orthologous genes" to "Through comparison of expression levels of 4,746 transcribed single-copy orthologous genes" (Main text, page 9, lines 200-201); from "For gene expression, there were critical biological differences among tissues (Pearson's  $r = 0.71$  and weighted average proportion variance = 0.42), followed by species (Pearson's  $r = 0.84$ , weighted average proportion variance = 0.16) and local adaptation (Pearson's  $r = 0.97$  and weighted average proportion variance = 0.019)" to "For gene expression, there were critical biological differences among tissues (Pearson's  $r = 0.67$  and weighted average proportion variance = 0.36), followed by species (Pearson's  $r = 0.75$ , weighted average proportion variance = 0.22) and local adaptation (Pearson's  $r = 0.95$  and weighted average proportion variance = 0.019)" (Main text, page 9, lines 206-210); from "We identified ~1,512 DEGs between 30 low-versus high-altitude pairs (225 DEGs in liver of pigs to 4,014 DEGs in kidney of sheep) (Table 1). Notably, among five pairs of vertebrate, the highly-diverged yak and cattle exhibited the highest number of DEG (~2,242) across six tissues. Among six tissues, the highly aerobic kidney exhibited the highest number of DEGs (~2,103) across five pairs of vertebrates." to "We identified ~1,423 DEGs between 30 low- versus high-altitude pairs (177 DEGs in muscle of chickens to 3,853 DEGs in kidney of sheep) (Table 1). Notably, among five pairs of vertebrate, the highly-diverged yak and cattle exhibited the highest number of DEG (~2,005) across six tissues. Among six tissues, the highly aerobic kidney exhibited the highest number of DEGs (~2,097) across five pairs of vertebrates" (Main text, page 11, lines 245-250).

The median of gene expression values (reflected by FPKM values) increased from 6.86 to 8.65, which corresponds to the increase of filtering cut-offs from 0.1 to 0.5 (Table R1 can be accessed from [RL\\_FiguresandTables.pdf](#) at: [https://www.dropbox.com/s/shgpb4784s409zw/RL\\_FiguresandTables.pdf?dl=0](https://www.dropbox.com/s/shgpb4784s409zw/RL_FiguresandTables.pdf?dl=0)).

#### Comment 2-5

b) On line 188, the authors use the term "high confidence single-copy orthologs" this is not defined. And is this homology based or expression based?

#### Response 2-5

Thank you for your valuable suggestions. We are sorry for our descriptive statement of approaches. We adopted the Ensemble pipeline that is more accurate than more

feasible OrthMCL method:

We applied the most recent Ensemble pipeline ([www.ensemble.org/info/genome/compara/homolog\\_method.html](http://www.ensemble.org/info/genome/compara/homolog_method.html)) to calculate 1:1 orthologues of five species. We downloaded the corresponding protein and CDS sequences of five species from Ensemble website with the exception of goat, whose protein and CDS sequences were downloaded from Goat Genome website. The sequences of an additional outgroup species zebrafish were also downloaded from Ensemble website. The longest protein sequence for each protein coding gene was kept for further analysis. Such protein sequences were concatenated to a single fasta file and makeblastdb function of NCBI blast+ version 2.2.28 [4] was applied to generate the reference file. The concatenated protein sequence fasta file was blasted against the reference file using blastp function of NCBI blast+: in effect, each gene of six species were blasted against each other (both within and between species), using parameters `-seg no -max_hsps_per_subject 1 -use_sw_tback -evalue 1e-10 -num_threads 1`. Blast e-values were converted to weights based on  $\text{MIN}(100, \text{ROUND}(-\text{LOG}_{10}(\text{evalue})/2))$ , and Hcluster\_sg (<http://sourceforge.net/p/treesoft/code/HEAD/tree/>) was utilized to cluster genes into families according to weights with parameters `-m 750 -w 0 -s 0.34`. Zebrafish was used as an outgroup species in this analysis by setting zebrafish genes to value 2 and non-zebrafish genes to value 1 in the category file, which was integrated into the analysis via `-C` option. Large clusters with more than 400 genes were recursively split into sub-clusters by QuickTree version 1.1 [5] until the largest sub-cluster contained less than 400 genes. In detail, multiple sequences of each large cluster were first aligned via Mafft version 7.149b [6] with parameter `-auto` and then converted to stockholm format by esl-reformat function in hmmer version 3.1b1 [7]. QuickTree were used to build unrooted tree and custom python scripts were utilized to find the branch that roughly split the tree into two parts of comparable nodes, by making sure one of the two parts contained the smallest possible number of nodes over half of the total number. This splitting process was repeated until the largest of the final sub-clusters had less than 400 genes. The split clusters were combined with the original clusters with less than 400 genes. Multiple alignment of protein sequences for each cluster was then generated by Mafft if there were over 200 genes, or by a mixture of four aligners of `mafftgins_msa`, `muscle_msa`, `kalign_msa` and `t_coffee_msa` consensified of M-coffee version 10.00.r1613 [8] if otherwise. For each aligned cluster, we back-translated the protein sequences to CDS and applied TreeBeST (<http://treesoft.sourceforge.net/treebest.shtml>) to build phylogenetic trees reconciled with an inputted species tree. Custom python scripts were utilized to retrieve one-to-one orthologues.

We also added the detailed method to the Supplementary Methods, hoping such information will help readers better understand our work.

Comment 2-6

2. Comparison of expression differences between high and low altitude animals and functional annotation analysis.

a) Supplemental Figure S3 shows that in some tissues there are large differences in mapping rate that are not reflected in the other altitude type. Did the authors check that mapping rate did not affect their differential expression calls? Also, please report the tissue type in this graph.

Response 2-6

As you suggested, we redrew the figures and compared the mapping ratios between low- and high-altitude populations for each vertebrate. Interestingly, we found that populations with a relatively lower mapping ratio of RNA-seq data had relatively higher genomic divergence from the reference genome (which was reflected by more SNPs based on whole-genome sequence data), and vice versa (Supplementary Fig. S3).

Thank you for pointing out that several tissues exhibited relatively lower mapping ratios. For example, hearts of high- and low-altitude pigs (Illumina HiSeq 2000 with 100-nt paired-end reads) and kidneys of low-altitude goats (Illumina HiSeq 2500 with 125-nt paired-end reads) (Supplementary Fig. S3) exhibited the lowest mapping ratios. This result indicated that the relatively lower mapping ratios may not be attributed to the idiosyncrasies of the different sequencing platforms.

We then considered that the discrepancies in mapping ratios might be attributable to bias from library construction, which can be effectively corrected during the

normalization steps implemented in cuffdiff [9]: to correct for library sizes (i.e., sequencing depths), FPKMs and fragment counts are scaled via the median of the geometric means of fragment counts across all libraries, as described by Anders and Huber [10].

#### Comment 2-7

b) In Additional File 2, a large table provided the GO/KEGG/InterPro terms and whether lists of genes with specific difference in high/low altitude expression are significantly enriched for that term. The authors should show the number of genes in the list for each comparison, or only show those with at least 5-10 genes in a list. Low representation in a pathway or term can be misleading for enrichment.

#### Response 2-7

Thank you for your valuable suggestions. We compared the similarities and differences of DE genes and their enriched categories between high altitude vertebrates and their low-altitude relatives within each tissue for each species (Supplementary Figs. S9-13, Additional Files 3-4). Then we retained gene lists with at least 5 genes, and updated all the relevant figure and tables accordingly.

#### Comment 2-8

c) More importantly, the authors do not indicate the background used for these analyses. It would be most appropriate to use the total number of genes expressed in each tissue for such analyses, so that the background reflects the genes that could possibly be shown to be differentially expressed, not the genome-wide background which is often the default.

#### Response 2-8

Thank you for your valuable comment. As previously reported [11-19], we used the annotated genes of whole-genome as the background for gene functional enrichment analysis in our initial submission. However, as you noted, it is more appropriate to use the genes expressed in each tissue as the background for gene functional enrichment analysis, which is more representative and could prevent the potential bias of over-representation of the tissue-specific expressed genes [20]. Based on your suggestion, we re-performed gene functional enrichment analysis by using ONLY the transcribed genes as the background, and found that the updated results were consistent with our initial results (Supplementary Figs. 12-13 and Additional Files 2, 4).

#### Reviewer 2:

#### Comment 3-1

First of all let me congratulate you and all authors for this piece of research. I have although some questions that I believe are important in order to improve your manuscript:

In Data Analysis:

#### Response 3-1

Thank you so much for your positive comments.

#### Comment 3-2

page 4, lines 85-88: may you specify how the data filtering was performed? which software did you use, or in case you have used in house developed scripts may you please provide them as supplemental information?

#### Response 3-2

Thank you so much for your questions. I used prinseq-0.20.4 [21], cutadapt-1.12 [22] and in house developed script to perform the filtering. The parameters used are 'prinseq-lite.pl -fastq R1.fastq -fastq2 R2.fastq -out\_format 3 -ns\_max\_p 10 -out\_good output -out\_bad null', and 'cutadapt -a AGATCGGAAGAGCACACGTCTGAACTCCAGTCAC --overlap=10 --error-rate=0.1 --discard-trimmed --paired-output tmp.2.fastq -o tmp.1.fastq R1\_1.fastq R2\_2.fastq', 'cutadapt -a AGATCGGAAGAGCGTCGTGTAGGGAAAGAGTGTAGATCTCGGTGGTCGCCGTAT CATT --overlap=10 --error-rate=0.1 --discard-trimmed --paired-output result\_1\_filteradapt.fastq -o result\_2\_filteradapt.fastq tmp.2.fastq tmp.1.fastq'

(Supplementary Methods).

Comment 3-3

page 4 line 93, may you specify the parameters used for the analysis performed with EnsemblComparaGeneTrees method?

Response 3-3

Thank you for your valuable suggestions. We applied the most recent Ensembl pipeline ([www.ensembl.org/info/genome/compara/homolog\\_method.html](http://www.ensembl.org/info/genome/compara/homolog_method.html)) to calculate 1:1 orthologues of five species. We downloaded the corresponding protein and CDS sequences of five species from Ensembl website with the exception of goat, whose protein and CDS sequences were downloaded from Goat Genome website. The sequences of an additional outgroup species zebrafish were also downloaded from Ensembl website. The longest protein sequence for each protein coding gene was kept for further analysis. Such protein sequences were concatenated to a single fasta file and makeblastdb function of NCBI blast+ version 2.2.28[4] was applied to generate the reference file. The merged protein sequence fasta file was blasted against the reference file using blastp function of NCBI blast+: in effect, each gene of six species were blasted against each other (both within and between species), using parameters `-seg no -max_hsps_per_subject 1 -use_sw_tback -evalue 1e-10 -num_threads 1`. Blast e-values were converted to weights based on  $\text{MIN}(100, \text{ROUND}(-\text{LOG}_{10}(\text{evalue})/2))$ , and Hcluster\_sg (<http://sourceforge.net/p/treesoft/code/HEAD/tree/>) was utilized to cluster genes into families according to weights with parameters `-m 750 -w 0 -s 0.34`. Zebrafish was used as an outgroup species in this analysis by setting zebrafish genes to value 2 and non-zebrafish genes to value 1 in the category file, which was integrated into the analysis via `-C` option. Large clusters with more than 400 genes were recursively split into sub-clusters by QuickTree version 1.1 [5] until the largest sub-cluster contained less than 400 genes. In detail, multiple sequences of each large cluster were first aligned via Mafft version 7.149b [6] with parameter `-auto` and then converted to stockholm format by `esl-reformat` function in hmmer version 3.1b1 [7]. QuickTree were used to build unrooted tree and custom python scripts were utilized to find the branch that roughly split the tree into two parts of comparable nodes, by making sure one of the two parts contained the smallest possible number of nodes over half of the total number. This splitting process was repeated until the largest of the final sub-clusters had less than 400 genes. The split clusters were combined with the original clusters with less than 400 genes. Multiple alignment of protein sequences for each cluster was then generated by Mafft if there were over 200 genes, or by a mixture of four aligners of `mafftgins_msa`, `muscle_msa`, `kalign_msa` and `t_coffee_msa` consensified by M-coffee version 10.00.r1613 [8] if otherwise. For each aligned cluster, we back-translated the protein sequences to CDS and applied TreeBeST (<http://treesoft.sourceforge.net/treebest.shtml>) to build phylogenetic trees reconciled with an inputted species tree. Custom python scripts were utilized to retrieve one-to-one orthologues (Supplementary Methods).

Comment 3-4

page 4- line 96, may you please detail the parameters used for the BWA alignment?

Response 3-4

Thank you for the valuable suggestions. The parameters are `'bwa mem -t 10 -k 32 -M'` (Supplementary Methods).

Comment 3-5

page 5- line 101- which were the parameters defined for GATK detection of SNPs and Indels? Parameters like Calling confidence and minimum read depth?

Response 3-5

Thank you for your valuable suggestions. `AddOrReplaceReadGroups` and `BuildBamIndex` function in Picard version 1.14 (<http://sourceforge.net/projects/picard/>) was applied to add read group information and index, separately. Indel realignment was performed using `RealignerTargetCreator` and `IndelRealigner` tools in GATK. We called variants by `HaplotypeCaller`, separated SNVs and Indels using `SelectVariants`, filtered SNVs with Fisher Strand values  $>60$  or Qual By Depth values  $<2$  or Mapping

Quality values<40 or Mapping Quality Rank Sum Test values<-12.5 or Read Position Rank Sum Test values<-8, and filtered Indels with Fisher Strand values>200 or Qual By Depth values<2 or Read Position Rank Sum Test values<-20 (Supplementary Methods).

Comment 3-6

page 5 line 108- which parameters were used for the TopHat alignment?

Response 3-6

Thank you for your valuable suggestions. The parameters we used are '--library-type fr-firststrand -p 4 --output-dir myoutputdir -G myspecies.gtf myspecies\_genomeindex read1.fq.gz read2.fq.gz' (Supplementary Methods).

Comment 3-7

In Findings:

I am missing analysis that I was expecting in a study of adaptation to altitude which generated so much WGS data. I suggest that you study genetic divergence by Fst or by Tajima's D and make identification of selection footprints. It would be great then to compare the genes being harbored in selective sweeps and the changes at transcriptomic level.

Response 3-7

We greatly appreciate your valuable comments.

At present, few studies have sufficiently characterized the direct relationship between genes embedded in selected regions and expression changes. Consequently, exploring the potential impact of positive selection on gene transcription is of great interest. As far as we know, only three vertebrates have publicly available whole-genome sequences for multiple individuals of both low-altitude populations (Pengxian chickens, Rongchang pigs, and Jersey cattle) and their high-altitude relatives (Tibetan chickens, Tibetan pigs, and yak) [23-26] (Table R2 can be accessed from RL\_FiguresandTables.pdf at: [https://www.dropbox.com/s/shgpb4784s409zw/RL\\_FiguresandTables.pdf?dl=0](https://www.dropbox.com/s/shgpb4784s409zw/RL_FiguresandTables.pdf?dl=0)).

To investigate the effects of positive selection on gene expression, we downloaded the above datasets and identified the genes embedded in selected regions (see Fig. R1) for high-altitude populations (Tibetan chickens, Tibetan pigs, and yak) against their low-altitude relatives (Pengxian chickens, Rongchang pigs, and Jersey cattle) (see Figs. R2-4) (Figs. R1-4 can be accessed from RL\_FiguresandTables.pdf at: [https://www.dropbox.com/s/shgpb4784s409zw/RL\\_FiguresandTables.pdf?dl=0](https://www.dropbox.com/s/shgpb4784s409zw/RL_FiguresandTables.pdf?dl=0)).

We found the genes embedded in selected regions exhibited highly comparable expression levels between the high-altitude populations and their low-altitude relatives within each tissue for each vertebrate, which was similar (P values of Wilcoxon rank sum test range from 0.120 to 0.939) to the genes outside selected regions (see Fig. R2).

We further observed expression levels of genes embedded in selected regions are highly comparable with the genes outside selected regions within each tissue for high-altitude population of each vertebrate (P values of Wilcoxon rank sum test range from 0.297 to 0.934) (see Fig. R3), this tendency also exists in their respective low altitude relatives (P values of Wilcoxon rank sum test range from 0.346 to 0.940) (see Fig. R4).

In this study, we did not observe the effects of positive selection on gene expression, which was most likely due to the distinct functional roles of variations with highly skewed frequency spectra. Generally, SNPs can be classified as coding (synonymous, missense, and nonsense) and non-coding. It is essential to perform further functional analyses to assess the impact of variations on gene expression; it is especially necessary to decipher the impact of non-coding variations that are located in regulatory regions (in particular, promoters, enhancers, and silencers) on gene expression.

Additionally, it is worth noting that our investigation is based on different individuals and had a small sample size; further large-scale experiments with proper design would be beneficial for answering this question.

Comment 3-8

page 10 lines 230-235: Did this happen in the low altitude chicken or only in one? its hard to see this in the figure

Response 3-8



Thank you for your thoughtful comment. As shown in the updated Fig. 4a and 4b (see Response 2–4), the Tibetan chickens and their low-altitude relatives formed a distinct cluster from the mammals. We revised this part of the manuscript to: “Notably, tissues of birds (chickens) formed a distinct cluster, rather than with their mammalian counterparts, which indicates that divergence in gene expression among these species started to surpass that between different tissues around when birds diverged from mammals (approximately 300 million years ago).” (Figs. 4a and 4b)

#### Comment 3-9

page 11 lines 251-259: The way these results are presented its hard to infer if the pathways affected by adaptation to altitude if these were the same between species or not. This is an important question that your results would enable to answer. I would suggest that a table per species should be made as well as venn diagrams that would lead us to understand which pathways were commonly affected or were different between species and if these were the same also at tissue level. I would like to see this part of the manuscript more enhanced, giving a larger value to the high value data that you have generated in your research.

#### Response 3-9

Thank you for your valuable suggestions, which are also commented by reviewer 1 (please see Response 2-3 as follows).

Thank you for your valuable suggestions. Based on your suggestions, we evaluated the amount of shared DE genes between the high- and low-altitude populations in each tissue among five vertebrates (Supplemental Figs. S9–10 and Additional File 3), and found that more closely related vertebrates shared more common DE genes (Supplemental Fig. S11). We also discovered that the enriched functional categories of DE genes substantially overlapped (Supplemental Figs. S12–13 and Additional File 4). We added Supplemental Figs. S9–13 and Additional Files 3–4 to the manuscript.

As shown in the newly added Supplemental Figs. S9-13 and Additional Files 3-4, expectedly, the more closely related vertebrates (Fig. 1) shared more DE genes (Supplementary Figs. S9–10 and Additional File 3). Compared with shared DE genes among mammals, especially between the two closely related members of Caprinae (goat and sheep), the birds (chickens) exhibited significantly fewer shared DE genes with mammals (Wilcoxon rank sum test,  $P < 0.0021$ ) (Supplementary Fig. S11). We also identified significantly enriched functional gene categories of DE genes (Chi-square test or Fisher’s exact test,  $P < 1.03 \times 10^{-4}$ ), which were shared among multiple pairwise comparisons (Supplementary Figs. S12–13 and Additional File 4), that were potentially related to the dramatic phenotypic changes shaped by high-altitude adaptation, such as response to hypoxia (typically, ‘oxidation reduction’, ‘heme binding’, ‘oxygen binding’, ‘oxygen transport’ and ‘oxygen transporter activity’), cardiovascular system (‘angiogenesis’ and ‘positive regulation of angiogenesis’), the efficiency of biomass production in the resource-poor highland (‘metabolic pathways’, ‘cholesterol biosynthetic process’ and ‘steroid metabolic process’) as well as immune response (‘responses of immune and defense’) (Additional file 2) (the statement has been added to the main text, page 11, line 251-267).

1. Brooks MJ, Rajasimha HK, Roger JE and Swaroop A. Next-generation sequencing facilitates quantitative analysis of wild-type and Nrl(-/-) retinal transcriptomes. *Mol Vis.* 2011;17:3034-54.

2. Hugo W, Shi H, Sun L, Piva M, Song C, Kong X, et al. Non-genomic and Immune Evolution of Melanoma Acquiring MAPKi Resistance. *Cell.* 2015;162(6):1271-85. doi:10.1016/j.cell.2015.07.061.

3. Mele M, Ferreira PG, Reverter F, DeLuca DS, Monlong J, Sammeth M, et al. Human genomics. The human transcriptome across tissues and individuals. *Science.* 2015;348(6235):660-5. doi:10.1126/science.aaa0355.

4. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al. BLAST+: architecture and applications. *BMC Bioinformatics.* 2009;10:421. doi:10.1186/1471-2105-10-421.

5. Howe K, Bateman A and Durbin R. QuickTree: building huge Neighbour-Joining trees of protein sequences. *Bioinformatics.* 2002;18(11):1546-7.

6. Katoh K and Standley DM. MAFFT multiple sequence alignment software version 7: improvements in performance and usability. *Mol Biol Evol.* 2013;30(4):772-80. doi:10.1093/molbev/mst010.

7. Finn RD, Clements J and Eddy SR. HMMER web server: interactive sequence

similarity searching. *Nucleic Acids Res.* 2011;39(Web Server issue):W29-37. doi:10.1093/nar/gkr367.

8.Wallace IM, O'Sullivan O, Higgins DG and Notredame C. M-Coffee: combining multiple sequence alignment methods with T-Coffee. *Nucleic Acids Res.* 2006;34(6):1692-9. doi:10.1093/nar/gkl091.

9.Trapnell C, Hendrickson DG, Sauvageau M, Goff L, Rinn JL and Pachter L. Differential analysis of gene regulation at transcript resolution with RNA-seq. *Nat Biotechnol.* 2013;31(1):46-53. doi:10.1038/nbt.2450.

10.Anders S and Huber W. Differential expression analysis for sequence count data. *Genome Biol.* 2010;11(10):R106. doi:10.1186/gb-2010-11-10-r106.

11.Yamaji M, Jishage M, Meyer C, Suryawanshi H, Der E, Yamaji M, et al. DND1 maintains germline stem cells via recruitment of the CCR4-NOT complex to target mRNAs. *Nature.* 2017;543(7646):568-72. doi:10.1038/nature21690.

12.Zhang Y, Sloan SA, Clarke LE, Caneda C, Plaza CA, Blumenthal PD, et al. Purification and characterization of progenitor and mature human astrocytes reveals transcriptional and functional differences with mouse. *Neuron.* 2016;89(1):37-53. doi:10.1016/j.neuron.2015.11.013.

13.Andor N, Graham TA, Jansen M, Xia LC, Aktipis CA, Petritsch C, et al. Pan-cancer analysis of the extent and consequences of intratumor heterogeneity. *Nat Med.* 2016;22(1):105-13. doi:10.1038/nm.3984.

14.Treutlein B, Lee QY, Camp JG, Mall M, Koh W, Shariati SA, et al. Dissecting direct reprogramming from fibroblast to neuron using single-cell RNA-seq. *Nature.* 2016;534(7607):391-5. doi:10.1038/nature18323.

15.Okin D and Medzhitov R. The effect of sustained inflammation on hepatic mevalonate pathway results in hyperglycemia. *Cell.* 2016;165(2):343-56. doi:10.1016/j.cell.2016.02.023.

16.Pimentel H, Parra M, Gee SL, Mohandas N, Pachter L and Conboy JG. A dynamic intron retention program enriched in RNA processing genes regulates gene expression during terminal erythropoiesis. *Nucleic Acids Res.* 2016;44(2):838-51. doi:10.1093/nar/gkv1168.

17.Mahat DB, Salamanca HH, Duarte FM, Danko CG and Lis JT. Mammalian heat shock response and mechanisms underlying its genome-wide transcriptional regulation. *Mol Cell.* 2016;62(1):63-78. doi:10.1016/j.molcel.2016.02.025.

18.Ortiz-Ramirez C, Hernandez-Coronado M, Thamm A, Catarino B, Wang M, Dolan L, et al. A transcriptome atlas of *Physcomitrella patens* provides insights into the evolution and development of land plants. *Mol Plant.* 2016;9(2):205-20. doi:10.1016/j.molp.2015.12.002.

19.Lacar B, Linker SB, Jaeger BN, Krishnaswami S, Barron J, Kelder M, et al. Nuclear RNA-seq of single neurons reveals molecular signatures of activation. *Nat Commun.* 2016;7:11022. doi:10.1038/ncomms11022.

20.Timmons JA, Szkop KJ and Gallagher IJ. Multiple sources of bias confound functional enrichment analysis of global -omics data. *Genome Biol.* 2015;16:186. doi:10.1186/s13059-015-0761-7.

21.Schmieder R and Edwards R. Quality control and preprocessing of metagenomic datasets. *Bioinformatics.* 2011;27(6):863-4. doi:10.1093/bioinformatics/btr026.

22.Martin M. Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnetJournal.* 2011;17(1):10-2.

23.Qiu Q, Wang L, Wang K, Yang Y, Ma T, Wang Z, et al. Yak whole-genome resequencing reveals domestication signatures and prehistoric population expansions. *Nat Commun.* 2015;6:10283. doi:10.1038/ncomms10283.

24.Daetwyler HD, Capitan A, Pausch H, Stothard P, van Binsbergen R, Brondum RF, et al. Whole-genome sequencing of 234 bulls facilitates mapping of monogenic and complex traits in cattle. *Nat Genet.* 2014;46(8):858-65. doi:10.1038/ng.3034.

25.Ai H, Fang X, Yang B, Huang Z, Chen H, Mao L, et al. Adaptation and possible ancient interspecies introgression in pigs identified by whole-genome sequencing. *Nat Genet.* 2015;47(3):217-25. doi:10.1038/ng.3199.

26.Li D, Che T, Chen B, Tian S, Zhou X, Zhang G, et al. Genomic data for 78 chickens from 14 populations. *GigaScience.* 2017;6(6):1-5. doi:10.1093/gigascience/gix026.

**Additional Information:**

**Question**

Are you submitting this manuscript to a

**Response**

No



special series or article collection?	
<p><b>Experimental design and statistics</b></p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>. 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><b>Resources</b></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 <a href="#">Research Resource Identifiers</a> (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p>	Yes
<p><b>Availability of data and materials</b></p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in <a href="#">publicly available repositories</a> (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 <a href="#">Minimum Standards Reporting Checklist</a>?</p>	Yes

---

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

3 Qianzi Tang<sup>1†</sup>, Yiren Gu<sup>2†\*</sup>, Xuming Zhou<sup>3†</sup>, Long Jin<sup>1</sup>, Jiuqiang Guan<sup>4</sup>, Rui Liu<sup>1</sup>, Jing Li<sup>1</sup>,  
4 Kereng Long<sup>1</sup>, Shilin Tian<sup>1</sup>, Tiandong Che<sup>1</sup>, Silu Hu<sup>1</sup>, Yan Liang<sup>2</sup>, Xuemei Yang<sup>2</sup>, Xuan  
5 Tao<sup>2</sup>, Zhijun Zhong<sup>2</sup>, Guosong Wang<sup>1,5</sup>, Xiaohui Chen<sup>2</sup>, Diyan Li<sup>1</sup>, Jideng Ma<sup>1</sup>, Xun  
6 Wang<sup>1</sup>, Miaomiao Mai<sup>1</sup>, An'an Jiang<sup>1</sup>, Xiaolin Luo<sup>4</sup>, Xuebin Lv<sup>2</sup>, Vadim N. Gladyshev<sup>3</sup>,  
7 Xuewei Li<sup>1\*</sup> and Mingzhou Li<sup>1\*</sup>

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

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

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

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

16 <sup>5</sup> 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](mailto:mingzhou.li@sicau.edu.cn), Xuewei Li:  
20 [xuewei.li@sicau.edu.cn](mailto:xuewei.li@sicau.edu.cn), Yiren Gu: [guyiren1128@163.com](mailto: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 were potentially involved in phenotypic divergence shaped by high-

---

49 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) [10], reads with mapping quality >  
99 0 were retained and potential PCR duplication cases were removed. For each  
100 individual, ~97.01% of reads were mapped to ~97.40% (at least 1 × depth  
101 coverage) or ~91.86% (at least 4 × depth coverage) of the reference genome  
102 assemblies (**Supplementary Table S2**). Single nucleotide variations (SNVs)  
103 and insertion and deletions (InDels) were further detected by following GATK's  
104 best practice (version 3.3-0) [11]. We substituted SNVs and InDels identified in  
105 our study in the coding DNA sequences (CDS) of the respective reference  
106 genomes. Single copy orthologues with substituted CDS of the five vertebrates  
107 were applied to Treebest [12] and generating the neighbor-joining tree (**Fig. 1b**).

#### 108 ***Analyses of gene expression***

109 High-quality RNA-seq reads were mapped to their respective reference  
110 genomes using Tophat (version 2.0.11) [13]. Cufflinks (version 2.2.1) [14] was  
111 applied to quantify gene expression and obtain FPKM expression values. We  
112 generated abundance files by applying Cuffquant (part of Cufflinks) to read  
113 mapping results. **Log<sub>2</sub>-transformed values of (FPKM + 1) for genes with >0.5**  
114 **FPKM in over 80% of the samples were used for subsequent analyses.**

115 Pearson's correlations were calculated across six samples from low- and  
116 high-altitudes populations within each group of specific tissue and animals;  
117 among pairwise comparisons of five animals within each of the six tissues; and  
118 among pairwise comparisons of six tissues within each of the five animals.  
119 Principal Variance Component Analysis (PVCA) was carried out using R  
120 package pvca [15]. Neighbor-joining expression-based trees were generated  
121 according to distance matrices composed of pairwise (1-Spearman's  
122 correlations) implemented in R package ape [16]. Reproducibility of branching

---

123 patterns was estimated by bootstrapping genes, that is, the single copy  
124 orthologues were randomly sampled with replacement 100 times. The fractions  
125 of replicate trees that share the branching patterns of the original tree  
126 constructed were marked by distinct node colors in the figure.

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

132 Genes were converted to human orthologs, and assessed by DAVID [17]  
133 webserver for functional enrichment in GO (Gene Ontology) terms consisting  
134 of molecular function (MF) and biological process (BP) as well as the KEGG  
135 pathways and InterPro databases (Benjamini adjusted p-values  $\leq 0.05$ ).

### 136 ***Analyses of alternative splicing***

137 Single-copy orthologous exons were identified by finding annotated exons that  
138 overlapped with the query exonic region in a multiple alignment of 99 vertebrate  
139 genomes including human genome (hg38) from the UCSC genome browser  
140 [18]. Exon groups with multiple overlapping exons in any species were  
141 excluded. Each internal exon in every annotated transcript was taken as an  
142 “cassette” exon. Each “cassett” alternative splicing (AS) is composed of three  
143 exons: C1, A and C2, where A is alternative exon, C1 the 5’ alternative exon,  
144 C2 the 3’ alternative exon. For each species and read length k, we generated  
145 all non-redundant constitutive and alternative junction sequences for the



---

146 following RNA-seq alignments. The junction sequences were constructed by  
147 retrieving k-8 bp from each of the two exons making up the junction, and when  
148 the exon length is smaller than k-8, the whole sequence of the exon is retrieved.  
149 This ensures that there is at least 8 bp overlap between the mapped reads and  
150 each of the two junction exons.

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

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

164 The “percent-spliced in” (PSI) values for each internal exon was defined as  
165  $PSI = 100 \times \text{average}(\#C1A, \#AC2) / (\#C1C2 + \text{average}(\#C1A, \#AC2))$ , here  
166 #C1A, #AC2 and #C1C2 are the normalized read counts for the associated  
167 junctions. Exons were taken as alternative in a sample if  $5 \leq PSI \leq 95$ . We also  
168 defined “high-confidence” PSI levels as those that meet the following criteria:

169  $\text{*max}(\text{min}(\#C1A, \#AC2), \#C1C2) \geq 5 \text{ AND } \text{min}(\#C1A, \#AC2) + \#C1C2 \geq 10$

170  $\text{*}|\log_2(\#C1A / \#AC2)| \leq 1 \text{ OR } \text{max}(\#C1A, \#AC2) < \#C1C2$

---

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

174

175

176

177

## Findings

### 178 *Data summary*

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

### 184 *Concordance in the tree topology based on nucleotide sequence* 185 *alignments and gene expression data*

186 Nucleotide alignments-based phylogenetic relationships of these high-altitude  
187 vertebrates and their low-altitude relatives matched the established  
188 morphological species groupings and the known history of population formation  
189 (**Fig. 1b**). The gene expression-based tree based 4,746 transcribed single-copy  
190 orthologous genes (66.61% of 7125) for each tissue showed a highly consistent  
191 topology to the nucleotide sequence alignment-based phylogeny (**Fig. 2,**  
192 **Supplementary Methods**) [9]: mammals were mainly divided into omnivore  
193 (pig) and ruminant (goat, sheep and yak/cattle); within the ruminant cluster, the  
194 two caprinae (goat and sheep) were closer to each other than the boviniae  
195 (yak/cattle). This observation lends supports the idea that gene expression

---

196 changes evolve together with genetic variation over evolutionary time, resulting  
197 in lower expression divergence between more closely species [19].

198 ***Distinctly transcriptomic characteristics between gene expression and***  
199 ***alternative splicing***

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

206 For gene expression, there were critical biological differences among tissues  
207 (Pearson's  $r = 0.67$  and weighted average proportion variance = 0.36), followed  
208 by species (Pearson's  $r = 0.75$ , weighted average proportion variance = 0.22)  
209 and local adaptation (Pearson's  $r = 0.95$  and weighted average proportion  
210 variance = 0.019) (**Fig. 3a** and **Supplementary Fig. S5**). By contrast, for  
211 alternative splicing, the differences among species (Pearson's  $r = 0.64$  and  
212 weighted average proportion variance = 0.30) were higher than among tissues  
213 (Pearson's  $r = 0.78$  and weighted average proportion variance = 0.075),  
214 followed by between high- and low-altitude animals (Pearson's  $r = 0.84$  and  
215 weighted average proportion variance = 0.021) (**Fig. 3b** and **Supplementary**  
216 **Figure S6**).

217 Unsupervised clustering (**Figs. 4a and 4c**) and principal components  
218 analysis (PCA) (**Figs. 4b and 4d** and **Supplementary Figs. S7 and S8**) both  
219 recapitulated the distinctly transcriptomic characteristics between gene  
220 expression and alternative splicing. Tissue-dominated clustering of gene

---

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

221 expression indicated that in general tissues possess conserved gene  
222 expression signatures and suggested that conserved gene expression  
223 differences underlie tissue identity in mammals. On the other hand, greater  
224 prominence of species-dominated clustering of alternative splicing suggested  
225 that exon splicing is more often affected by species-specific changes in *cis*-  
226 regulatory elements and/or *trans*-acting factors than gene expression [20, 21].

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

### 237 ***Gene expression plasticity to a high-altitude environment***

238 To exclude the impact of prominence of tissues-dominated clustering of gene  
239 expression, so as to comprehensively present transcriptomic differences  
240 involved in high-altitude response based on whole annotated genes of their  
241 respective genome assembly instead of the single-copy orthologous, we  
242 measured the pairwise difference of gene expression between the high-altitude  
243 populations and their low-altitude relatives within each tissue for each  
244 vertebrate.

---

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

251 Expectedly, the more closely related vertebrates (**Fig. 1**) shared more DE  
252 genes (**Supplementary Figs. S9–10 and Additional File 3**). Compared with  
253 shared DE genes among mammals, especially between the two closely related  
254 members of Caprinae (goat and sheep), the birds (chickens) exhibited  
255 significantly fewer shared DE genes with mammals (Wilcoxon rank sum test,  
256  $P < 0.0021$ ) (**Supplementary Fig. S11**). We also identified significantly enriched  
257 functional gene categories of DE genes (Chi-square test or Fisher's exact test,  
258  $P < 1.03 \times 10^{-4}$ ), which were shared among multiple pairwise comparisons  
259 (**Supplementary Figs. S12–13 and Additional File 4**), that were potentially  
260 related to the dramatic phenotypic changes shaped by high-altitude adaptation,  
261 such as response to hypoxia (typically, 'oxidation reduction', 'heme binding',  
262 'oxygen binding', 'oxygen transport' and 'oxygen transporter activity'),  
263 cardiovascular system ('angiogenesis' and 'positive regulation of  
264 angiogenesis'), the efficiency of biomass production in the resource-poor  
265 highland ('metabolic pathways', 'cholesterol biosynthetic process' and 'steroid  
266 metabolic process') as well as immune response ('responses of immune and  
267 defense') (**Additional file 2**).

## 268 **Conclusions**

---

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

277 **Availability of supporting data**

278 The RNA-seq data for 180 samples was deposited in the NCBI Gene  
279 Expression Omnibus (GEO) under accession numbers GSE93855, GSE77020  
280 and GSE66242. The re-sequencing data for 7 individuals was deposited in the  
281 NCBI-sequence read archive (SRA) under accession number SRP096151.

282 All supplementary figures and tables are provided in Additional file.

283 **Reviewer links:**

284 GSE93855:

285 <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?token=irqtigkqvtatnqt&acc=G>

286 [SE93855](#)

287 GSE77020:

288 <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?token=kpolsqsothybrcv&acc=>

289 [GSE77020](#) (GSM1617847-GSM1617849 and GSM2042608-GSM2042610 are  
290 duplicates and represent the same samples)

291 GSE66242:

---

292 [https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?token=absxuuywtfyhncx&acc](https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?token=absxuuywtfyhncx&acc=293 <u>GSE66242</u>)  
293 [=GSE66242](#) (9 goat samples derived from individuals sampled at 2000m  
294 altitude were not included in this study)

295

296 **Ethics statement**

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

306 **Consent for publication**

307 Not applicable.

308 **Competing interests**

309 The authors declare that they have no competing interests.

310 **Funding**

311 This work was supported by grants from the National High Technology  
312 Research and Development Program of China (863 Program) (2013AA102502),  
313 the National Natural Science Foundation of China (31402046, 31522055,  
314 31601918, 31530073 and 31472081), the Science & Technology Support



---

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 Program of Sichuan (2016NYZ0042), the Youth Science Fund of Sichuan  
316 (2017JQ0011), the China Postdoctoral Science Foundation (2015M572486),  
317 China Agriculture Research System (CARS-36), the Program for Innovative  
318 Research Team of Sichuan Province (2015TD0012), the Program for Pig  
319 Industry Technology System Innovation Team of Sichuan Province (SCCXTD-  
320 005), the Project of Sichuan Education Department (15ZA0008, 15ZA0003,  
321 16ZA0025 and 16ZB0037), the National Program for Support of Top-notch  
322 Young Professionals and the Young Scholars of the Yangtze River.

323 **Authors' contributions**

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

329

330

331

332 **References**

333

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

- 
- 1 340 3. Li M, Chen L, Tian S, Lin Y, Tang Q, Zhou X, et al. Comprehensive variation  
2 341 discovery and recovery of missing sequence in the pig genome using multiple de  
3  
4 342 novo assemblies. *Genome Res.* 2017 May;27(5):865-874.  
5  
6 343 doi:10.1101/gr.207456.116.  
7  
8 344 4. International Chicken Genome Sequencing Consortium. Sequence and  
9  
10 345 comparative analysis of the chicken genome provide unique perspectives on  
11  
12 346 vertebrate evolution. *Nature.* 2004;432(7018):695-716. doi:10.1038/nature03154.  
13  
14 347 5. Groenen MA, Archibald AL, Uenishi H, Tuggle CK, Takeuchi Y, Rothschild MF, et  
15  
16 348 al. Analyses of pig genomes provide insight into porcine demography and evolution.  
17  
18 349 *Nature.* 2012;491(7424):393-8. doi:10.1038/nature11622.  
19  
20 350 6. Bovine Genome Sequencing and Analysis Consortium, Elsik CG, Tellam RL,  
21  
22 351 Worley KC, Gibbs RA, et al. The genome sequence of taurine cattle: a window to  
23  
24 352 ruminant biology and evolution. *Science.* 2009;324(5926):522-8.  
25  
26 353 doi:10.1126/science.1169588.  
27  
28 354 7. Dong Y, Xie M, Jiang Y, Xiao N, Du X, Zhang W, et al. Sequencing and automated  
29  
30 355 whole-genome optical mapping of the genome of a domestic goat (*Capra hircus*).  
31  
32 356 *Nat Biotechnol.* 2013;31(2):135-41. doi:10.1038/nbt.2478.  
33  
34 357 8. Jiang Y, Xie M, Chen W, Talbot R, Maddox JF, Faraut T, et al. The sheep genome  
35  
36 358 illuminates biology of the rumen and lipid metabolism. *Science.* 2014;344  
37  
38 359 (6188):1168-73. doi:10.1126/science.1252806.  
39  
40 360 9. Vilella AJ, Severin J, Ureta-Vidal A, Heng L, Durbin R and Birney E.  
41  
42 361 EnsemblCompara GeneTrees: Complete, duplication-aware phylogenetic trees in  
43  
44 362 vertebrates. *Genome Res.* 2009;19(2):327-35. doi:10.1101/gr.073585.107.  
45  
46 363 10. Li H and Durbin R. Fast and accurate long-read alignment with Burrows-Wheeler  
47  
48 364 transform. *Bioinformatics.* 2010;26(5):589-95. doi:10.1093/bioinformatics/btp698.  
49  
50 365 11. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernysky A, et al.  
51  
52 366 The Genome Analysis Toolkit: a MapReduce framework for analyzing next-  
53  
54 367 generation DNA sequencing data. *Genome Res.* 2010;20(9):1297-303.  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

---

368 doi:10.1101/gr.107524.110.

369 12. Li H, Coghlan A, Ruan J, Coin LJ, Heriche JK, Osmotherly L, et al. TreeFam: a  
370 curated database of phylogenetic trees of animal gene families. *Nucleic Acids*  
371 *Res.* 2006;34(Database issue):D572-80. doi:10.1093/nar/gkj118.

372 13. Trapnell C, Pachter L and Salzberg SL. TopHat: discovering splice junctions with  
373 RNA-Seq. *Bioinformatics.* 2009;25(9):1105-11. doi:10.1093/bioinformatics/btp120.

374 14. Trapnell C, Williams BA, Pertea G, Mortazavi A, Kwan G, van Baren MJ, et al.  
375 Transcript assembly and quantification by RNA-Seq reveals unannotated  
376 transcripts and isoform switching during cell differentiation. *Nat Biotechnol.*  
377 2010;28(5):511-5. doi:10.1038/nbt.1621.

378 15. The `pvca` R package.  
379 <https://bioconductor.org/packages/release/bioc/html/pvca.html>. Accessed Feb 16  
380 2017.

381 16. Paradis E, Claude J and Strimmer K. APE: Analyses of phylogenetics and  
382 evolution in R language. *Bioinformatics.* 2004;20(2):289-90.

383 17. Dennis G, Jr., Sherman BT, Hosack DA, Yang J, Gao W, Lane HC, et al. DAVID:  
384 Database for annotation, visualization, and integrated discovery. *Genome Biol.*  
385 2003;4(5):P3.

386 18. Kent WJ, Sugnet CW, Furey TS, Roskin KM, Pringle TH, Zahler AM, et al. The  
387 human genome browser at UCSC. *Genome Res.* 2002;12(6):996-1006.  
388 doi:10.1101/gr.229102.

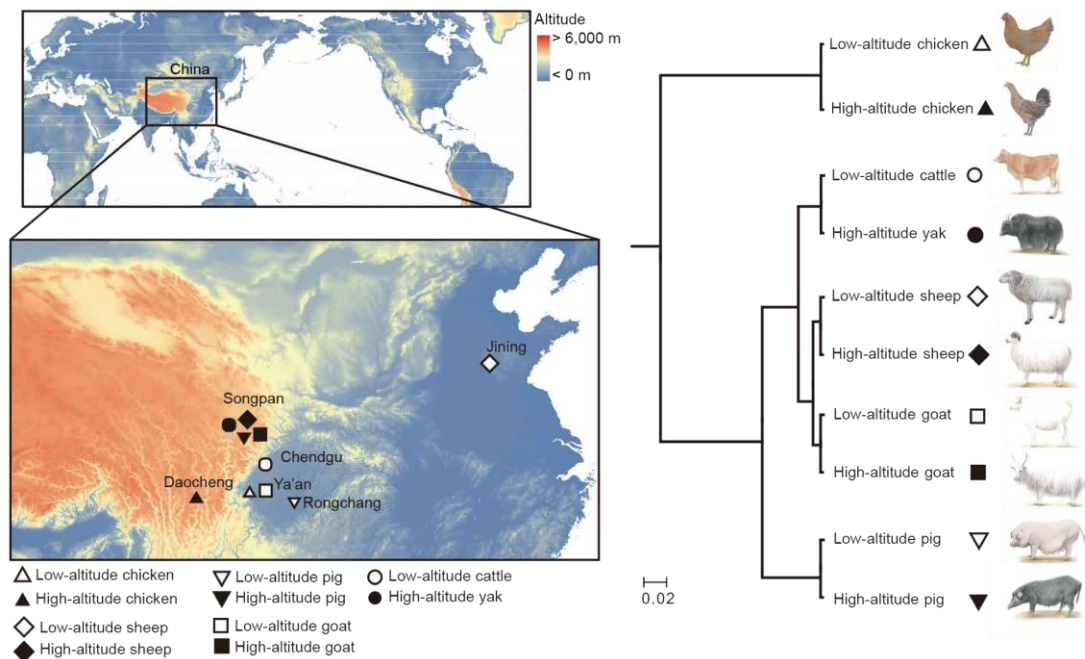
389 19. Brawand D, Soumillon M, Neacsulea A, Julien P, Csardi G, Harrigan P, et al. The  
390 evolution of gene expression levels in mammalian organs. *Nature.* 2011;478  
391 (7369):343-8. doi:10.1038/nature10532.

392 20. Merkin J, Russell C, Chen P and Burge CB. Evolutionary dynamics of gene and  
393 isoform regulation in Mammalian tissues. *Science.* 2012;338(6114):1593-9.  
394 doi:10.1126/science.1228186.

395 21. Barbosa-Morais NL, Irimia M, Pan Q, Xiong HY, Gueroussov S, Lee LJ, et al. The

396 evolutionary landscape of alternative splicing in vertebrate species. *Science*.  
 397 2012;338(6114):1587-93. doi:10.1126/science.1230612.  
 398 22. Mele M, Ferreira PG, Reverter F, DeLuca DS, Monlong J, Sammeth M, et al.  
 399 Human genomics. The human transcriptome across tissues and individuals.  
 400 *Science*. 2015;348(6235):660-5. doi:10.1126/science.aaa0355.

406 **Figures 1-4**



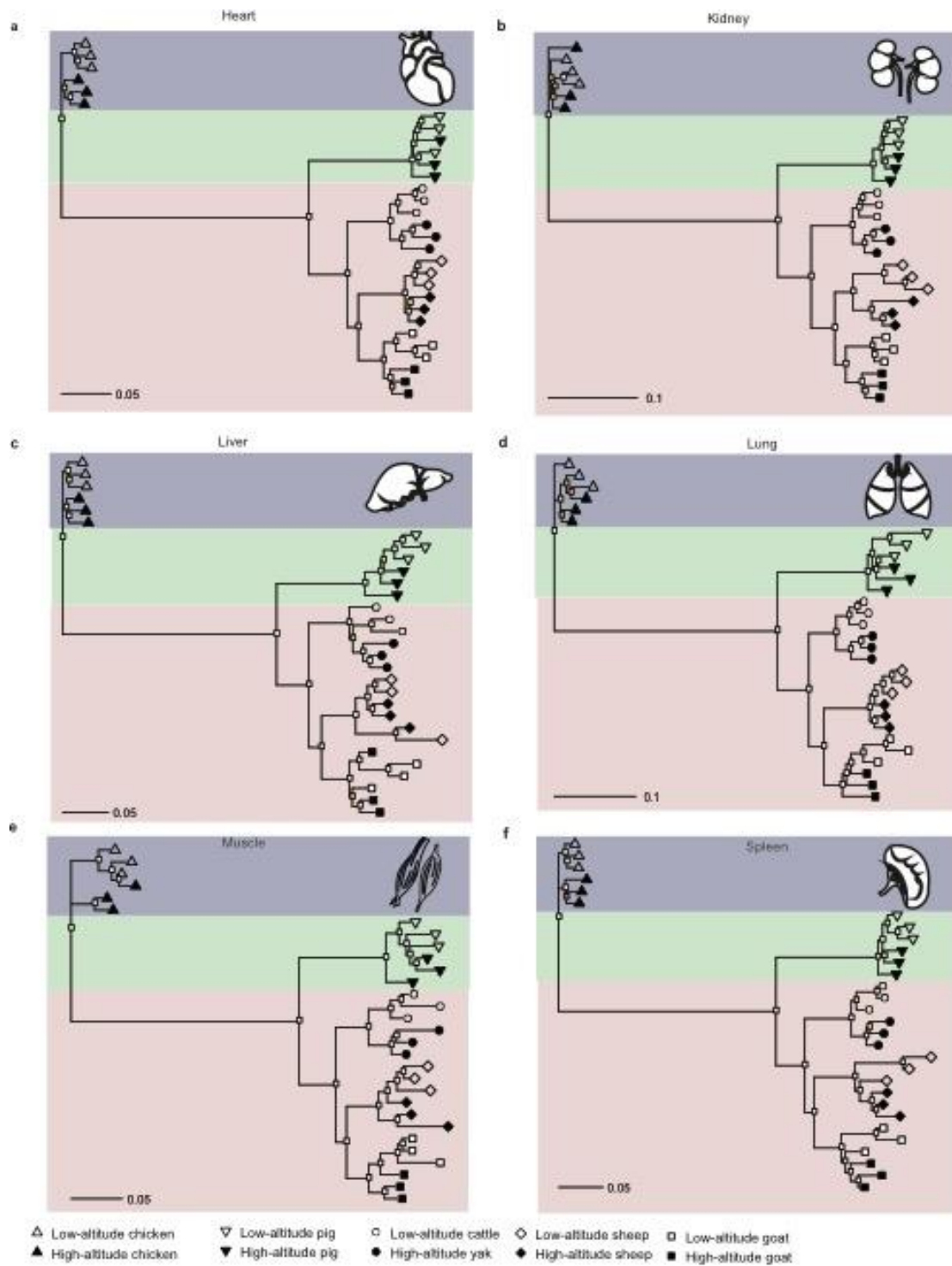
410 **Figure 1. Sampling locations and nucleotide alignment-based tree.**

411 **(a) Geographic locations of the studied animals.**

---

1  
2 412 (b) A neighbour-joining tree constructed based on concatenated coding sequences of  
3 413 single-copy orthologues substituted by SNVs and InDels detected in each animal.  
4 414 We downloaded and extracted the unassembled reads from short-insert (500 bp) libraries  
5  
6 415 of a single yak [1], a Tibetan pig [2] and a Rongchang pig [3] that were used for *de novo*  
7  
8 416 assemblies to roughly 10 × depth coverage. We also randomly selected an individual of  
9  
10 417 the cattle, low- and high-altitude chicken, goat and sheep and sequenced the whole  
11  
12 418 genomes at ~10 × depth coverage.  
13

14 419  
15  
16 420  
17  
18 421  
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



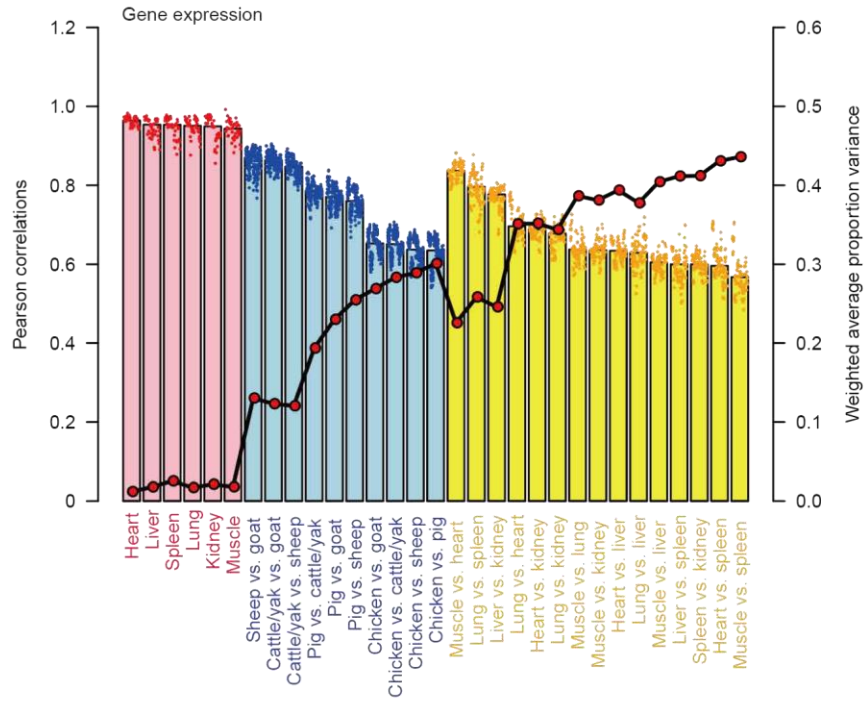
422

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

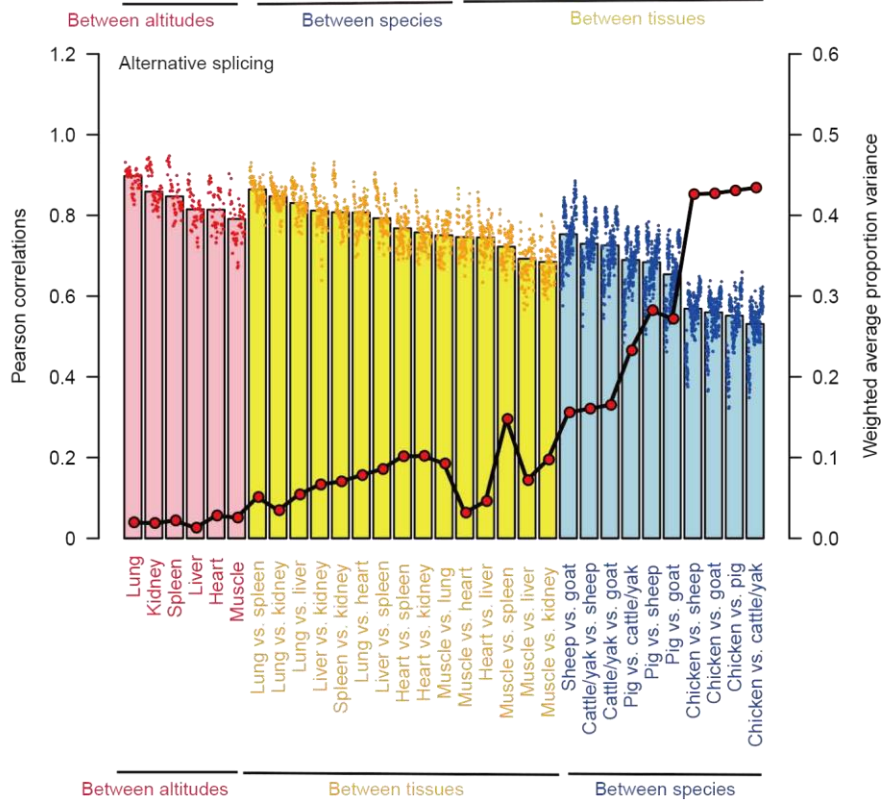
424 Neighbour-joining expression tree constructed based on (1-Spearman correlation)  
 425 distances in six tissues. We performed 100 bootstraps by randomly sampling the single  
 426 copy orthologues with replacement. Bootstrap values (fractions of replicate trees that have  
 427 the branching pattern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

429 less than 50% bootstraps, while orange, yellow and white colors indicate support between  
 430 50% and 70%, between 70% and 90% and more than 90%, respectively.

a



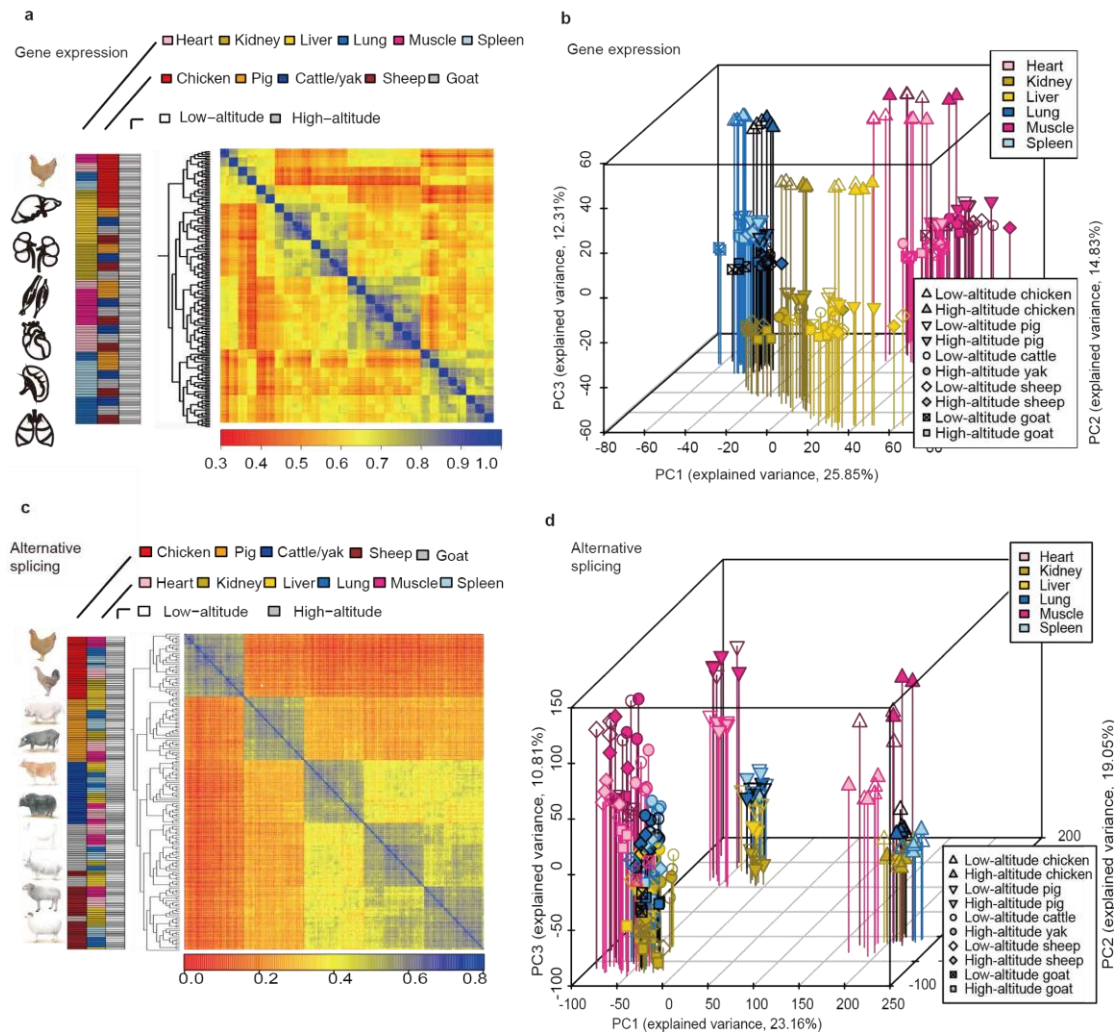
b



431  
 432 **Figure 3. Comparison of variations between altitudes, species and tissues revealed**  
 433 **by (a) gene expression and (b) alternative splicing pattern.**



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



438

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

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

---

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

444 expression levels and **(d)** the alternative splicing. The proportion of the variance

445 explained by the principal components is indicated in parentheses.

446

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

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

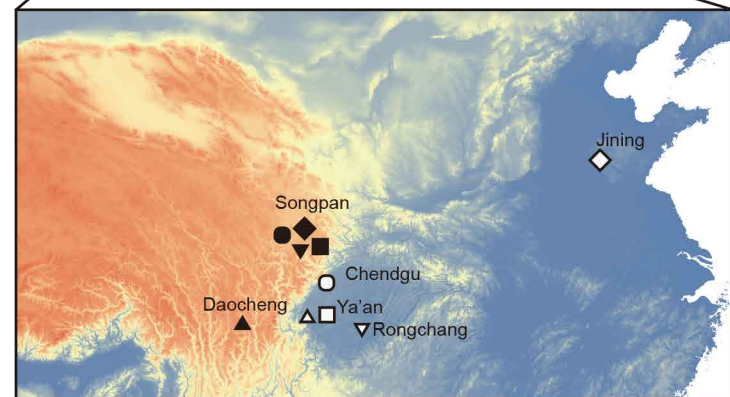
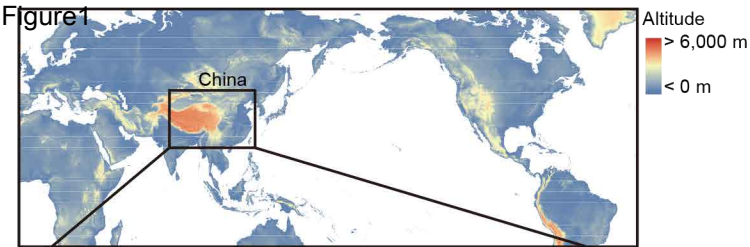
Percentage of the DGEs compared with the total number of annotated protein coding genes in their respective reference genomes are listed in parenthesis. There are 15495, 21594, 19981, 22131, 20908 annotated protein coding genes in reference genomes of Chicken (Galgal4) [4], pig (Suscrofa 10.2) [5], cattle (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**

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

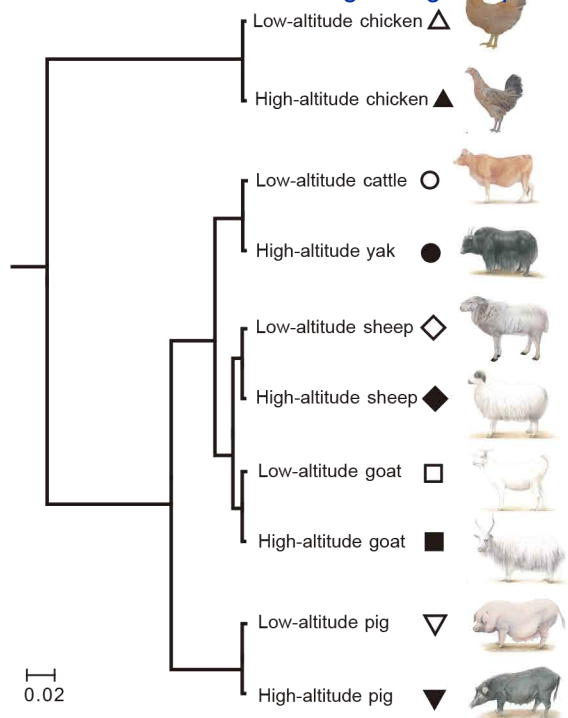
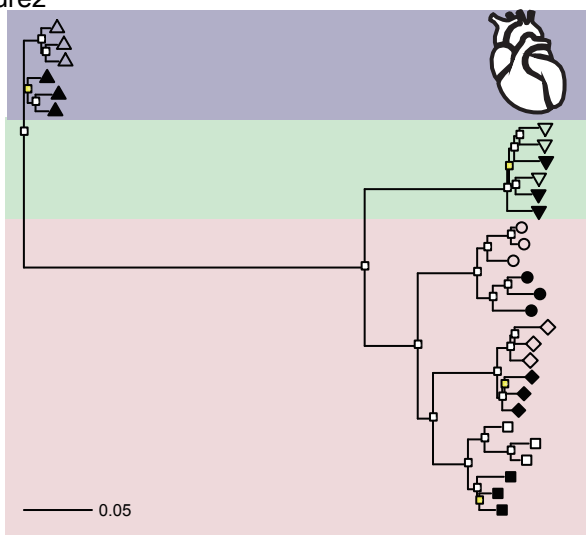
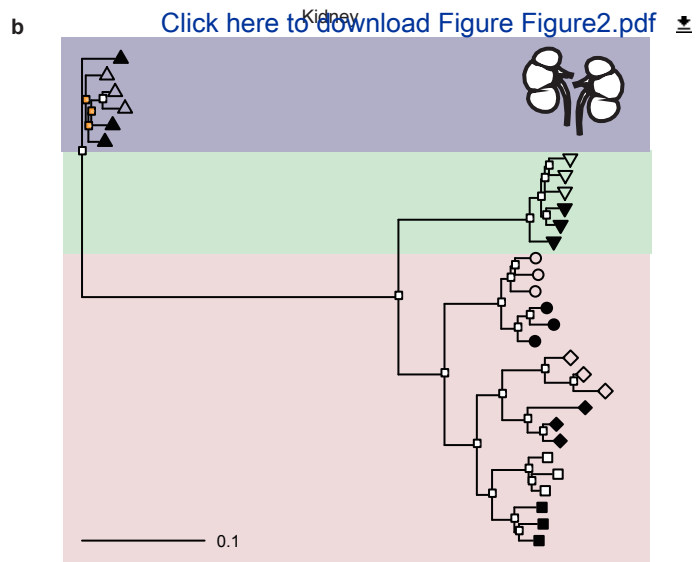


Figure 2

Heart



Kidney

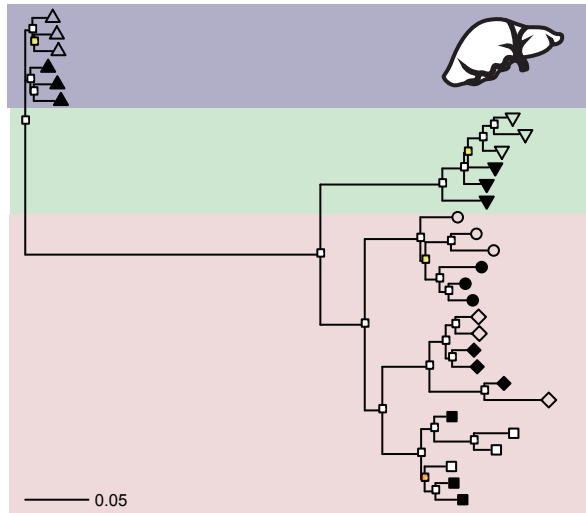


[Click here to download Figure 2.pdf](#)



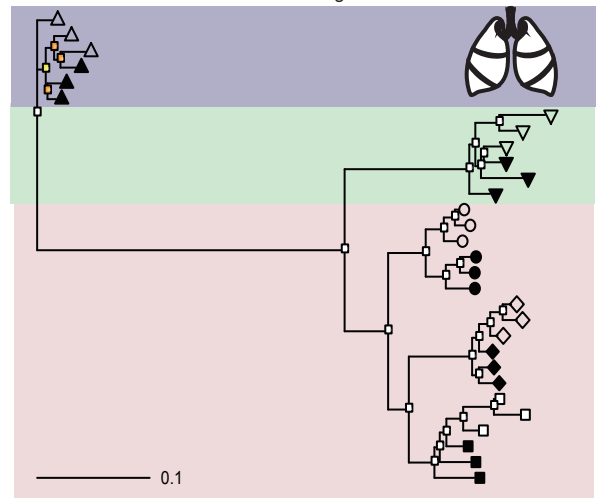
c

Liver



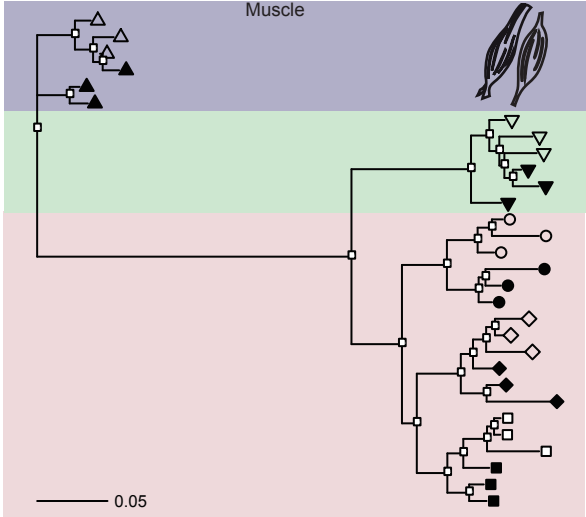
d

Lung



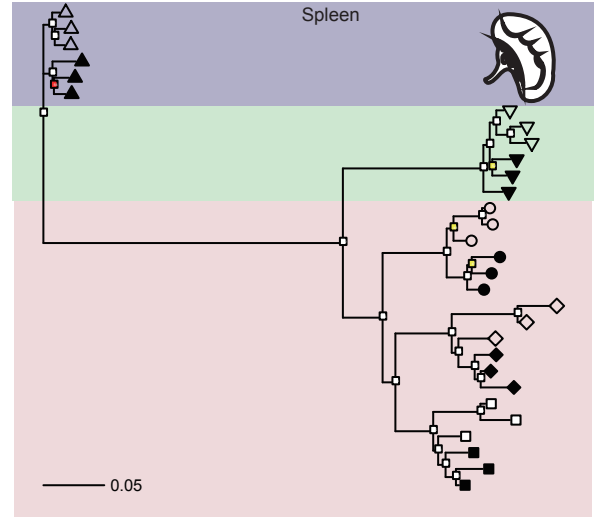
e

Muscle



f

Spleen



△ Low-altitude chicken  
▲ High-altitude chicken

▽ Low-altitude pig  
▼ High-altitude pig

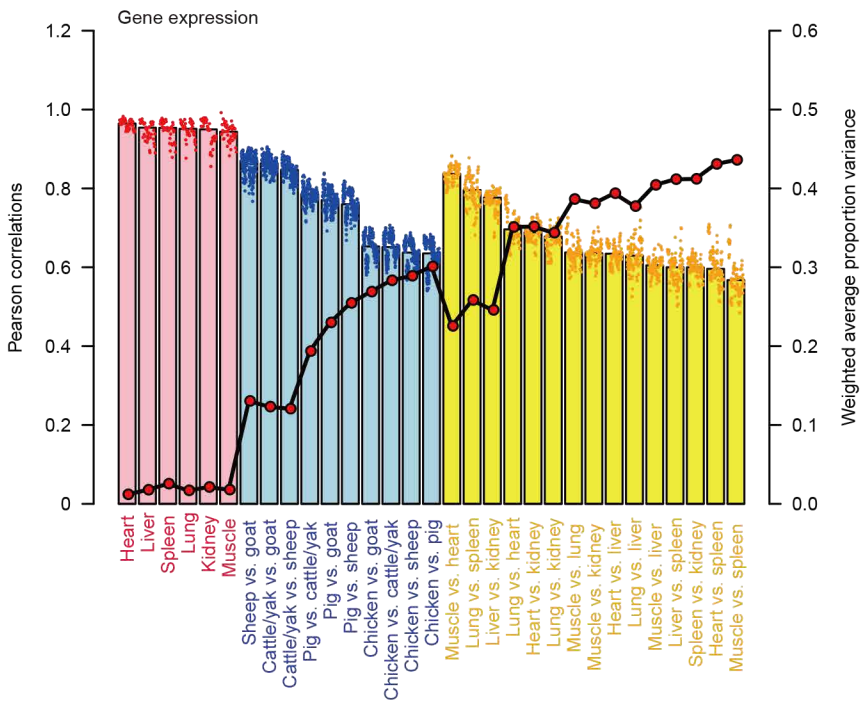
○ Low-altitude cattle  
● High-altitude yak

◇ Low-altitude sheep  
◆ High-altitude sheep

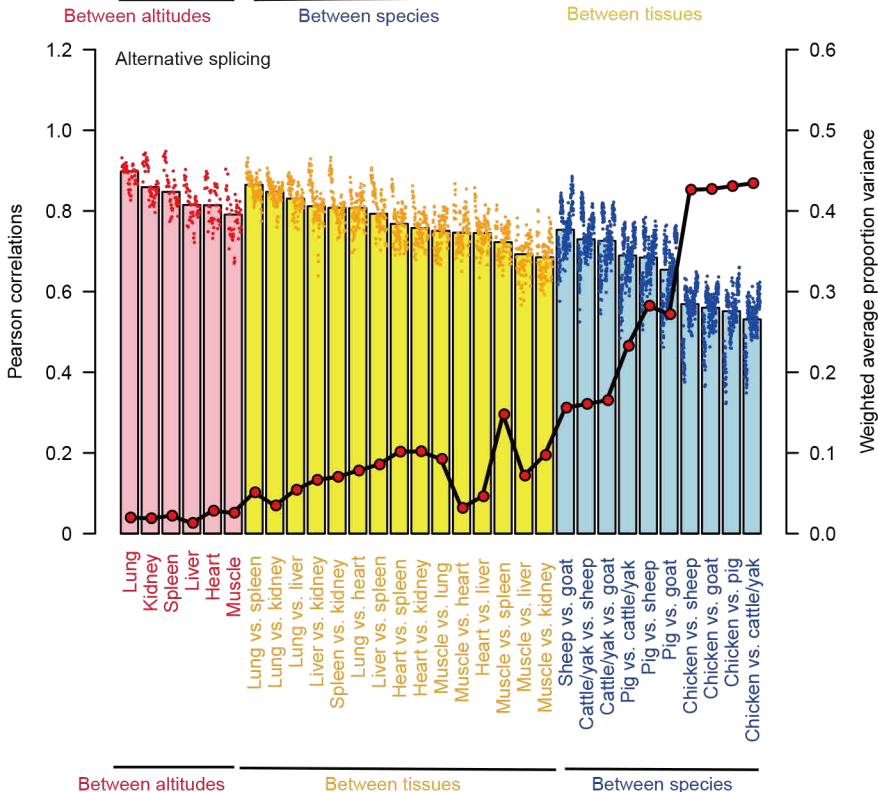
□ Low-altitude goat  
■ High-altitude goat

Figure3

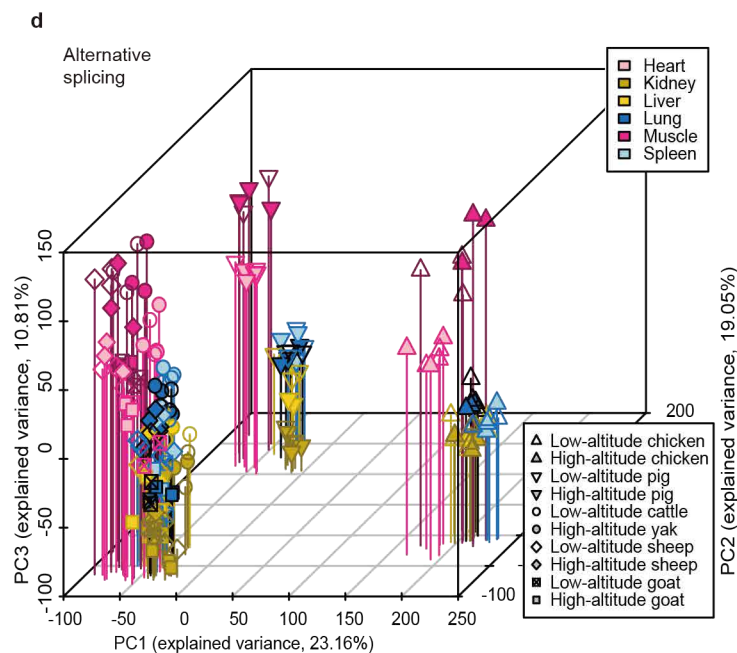
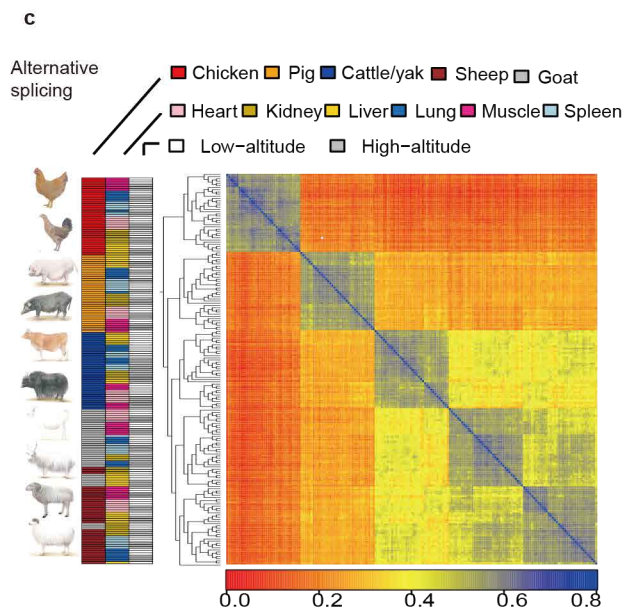
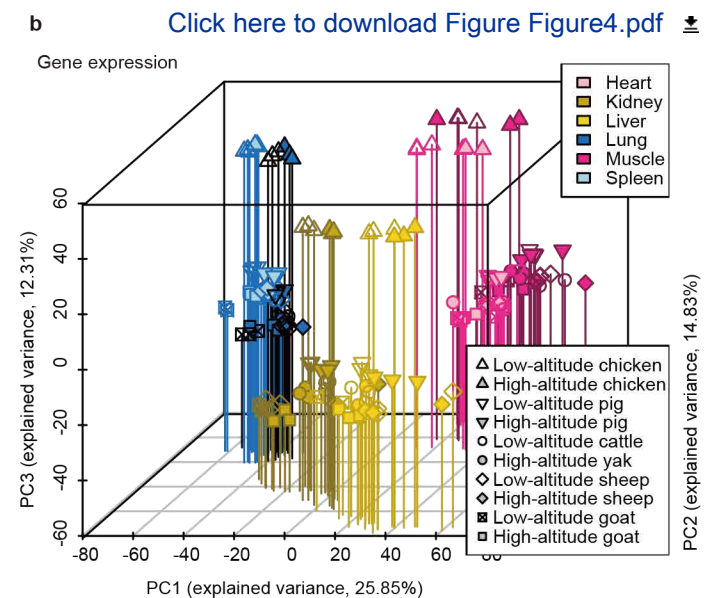
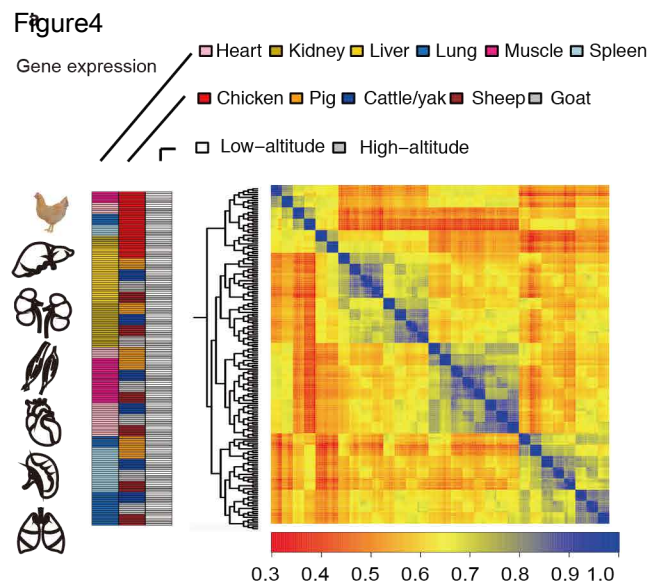
[Click here to download Figure Figure3.pdf](#)



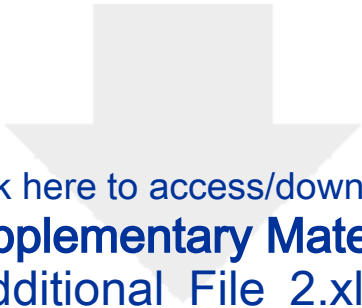
b






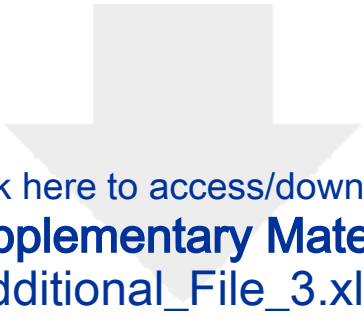







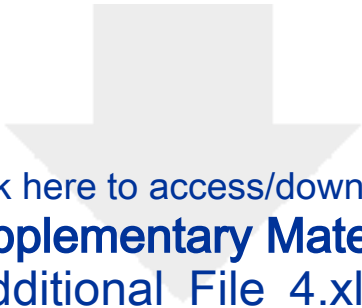
Click here to access/download  
**Supplementary Material**  
Additional\_File\_2.xlsx






Click here to access/download  
**Supplementary Material**  
Additional\_File\_3.xlsx





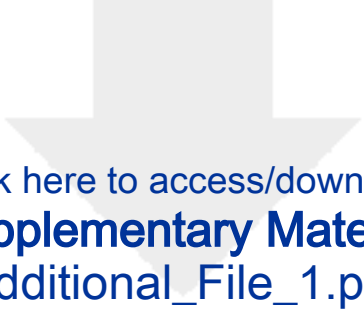
Click here to access/download  
**Supplementary Material**  
Additional\_File\_4.xlsx





Click here to access/download  
**Supplementary Material**  
RL\_FiguresandTables.pdf





Click here to access/download  
**Supplementary Material**  
Additional\_File\_1.pdf



***GigaScience***

em@editorialmanager.com

Dear Dr. Hans Zauner,

We are delighted to be informed of the positive responses from you. Thank you for your consideration of our manuscript for publication in *GigaScience*. We sincerely appreciate the thoughtful and constructive comments from the two reviewers Dr. Christopher Tuggle and Dr. Andreia Amaral, and your assistance in improving the manuscript. We have gone through your as well as the reviewers' comments in detail and believe that we have fully addressed these questions and concerns. We substantially improved our manuscript, and added 5 supplementary figures comprising 30 panels and 2 additional files. Below we provide our point-to-point responses, and hope that they are satisfactory.

We look forward to hearing a positive response from you.

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:

*Reviewer 1 points out that you should share all of your gene lists - I should clarify this point, as it may be confusing without explanation: at the time of writing their report, the reviewer did not have access to your additional files, and in fact you already provide a lot of this material, as the reviewer also has confirmed in further correspondence, after inspecting the files. However, please make sure the gene lists and other supporting data are in fact complete.*

#### Response 1-1:

Thank you for your reminder. According to the submission guidelines of *GigaScience*, we uploaded the complete gene lists with normalized expression values to the *GigaScience* temporary FTP server.

#### Comment 1-2:

*Please also add the additional information on your methods and analyses the reviewers are asking for. Reviewer 2 recommends some additional analyses, that I feel will be a useful addition.*

#### Response 1-2:

Thank you for your comments. Based on reviewer 2's valuable recommendations, we carried out two additionally explorative analyses.

First, to investigate the similarities of the gene expression pairwise differences between the high- and low-altitude populations, we identified

shared differentially expressed (DE) genes and common functional categories enriched by DE genes in the pairwise comparisons of each tissue for each vertebrate. We added **Supplementary Figs. 9–13** and **Additional Files 3–4** in the revised manuscript (see **Responses 2–3** and **3–9**, respectively, for details).

Second, to explore the potential impact of positive selection on gene transcription, we identified the genes embedded in the selected regions based on publicly available whole-genome sequence data of three vertebrates (i.e., chickens, pigs, and cattle that live at low altitudes) and their high-altitude relatives, and compared these genes and the changes at the transcriptomic level (see **Response 3–7**, **Figs. R1–4**, and **Table R2** for details, **Figs. R1–4**, and **Table R2** can be accessed from RL\_FiguresandTables.pdf at: [https://www.dropbox.com/s/shgpb4784s409zw/RL\\_FiguresandTables.pdf?dl=0](https://www.dropbox.com/s/shgpb4784s409zw/RL_FiguresandTables.pdf?dl=0)).

**Comment 1-3:**

*Please also make sure you follow all of the MNSEQE standards of reporting:*

*[http://www.fged.org/site\\_media/pdf/MINSEQE\\_1.0.pdf](http://www.fged.org/site_media/pdf/MINSEQE_1.0.pdf)*

**Response 1-3:**

Thank you for your kind reminder. According to your notification, we checked the data styles to completely conform to the MNSEQE standards of reporting.