

Genomic data for 78 chickens from 14 populations

1 Diyan Li^{1†}, Tiandong Che^{1†}, Binlong Chen^{1†}, Shilin Tian^{1,2†}, Xuming Zhou^{3†}, Guolong Zhang^{4†},
2 Miao Li¹, Uma Gaur¹, Yan Li¹, Majing Luo⁵, Long Zhang¹, Zhongxian Xu¹, Xiaoling Zhao¹,
3 Huadong Yin¹, Yan Wang¹, Long Jin¹, Qianzi Tang¹, Huailiang Xu¹, Mingyao Yang¹, Rongjia
4 Zhou⁵, Ruiqiang Li², Qing Zhu¹ and Mingzhou Li¹

5
6
7
8 ¹ Institute of Animal Genetics and Breeding, College of Animal Science and Technology,
9 Sichuan Agricultural University, Chengdu, China

10 ² Novogene Bioinformatics Institute, Beijing, China

11 ³ Division of Genetics, Department of Medicine, Brigham and Women's Hospital, Harvard
12 Medical School, Boston, USA

13 ⁴ Department of Animal Science, Oklahoma State University, Stillwater, Oklahoma, USA

14 ⁵ Hubei Key Laboratory of Cell Homeostasis, Laboratory of Molecular and Developmental
15 Genetics, College of Life Sciences, Wuhan University, Wuhan, China

16 [†] These authors contributed equally to this work.

17 **Correspondence:** zhuqingsicau@163.com; mingzhou.li@sicau.edu.cn.

Abstract

18 **Background:** Since the domestication of the red jungle fowls (*Gallus gallus*) (dating back to
19 ~10,000 B.P.) in Asia, domestic chickens (*Gallus gallus domesticus*) have been subjected to the
20 combined effects of natural selection and human-driven artificial selection; this has resulted in
21 marked phenotypic diversity in a number of traits, including behavior, body composition, egg
22 production and skin color. Population genomic variations through diversifying selection have
23 not been fully investigated.

24 **Findings:** The whole genomes of 78 domestic chickens were sequenced to an average of 18-
25 fold coverage for each bird. By combining this data with publicly available genomes of 5 wild
26 red jungle fowls and 8 Xishuangbanna game fowls, we conducted a comprehensive
27 comparative genomics analysis of 91 chickens from 17 populations. After aligning ~21.30
28 gigabases (Gb) of high quality data from each individual to the reference chicken genome, we
29 identified ~6.44 million (M) SNPs for each population. These SNPs included 1.10 M novel
30 SNPs in 17 populations that were absent in the current chicken dbSNP (Build 145) entries.

31 **Conclusions:** The current data is important for population genetics and further studies in
32 chicken, and will serve as a valuable resource for investigating diversifying selection and
33 candidate genes for selective breeding in chicken.

34 **Keywords:** Chicken, Genetic diversity, Population genomics, Whole-genome resequencing

Data description

Genome sequencing and sequence filtering

35 The 78 blood samples (36 Tibetan fowls from the Qinghai-Tibet Plateau and 42 domestic
36 fowls from Szechwan Basin) (Figure 1) were collected from the wing vein. The animal handling
37 experiments were approved by the Institutional Animal Care and Use Committee of Sichuan
38 Agricultural University under permit number YCS-B20100804. Genomic DNA was extracted
39 from these samples following standard procedures. In total, we generated ~1.69 trillion bases
40 of resequencing data of the whole genomes from 78 birds (18.03-fold coverage for each
41
42
43
44

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 individual) on the Illumina Hiseq 2500 platform (Additional file 1:Table S1). In addition,
46 previously published genome sequence data from 5 red jungle fowls (RJF) and 8
47 Xishuangbanna game fowls (~16.6-fold coverage for each individual) were downloaded and
48 analyzed (GenBank accession number PRJNA241474) (Figure 1).

49 We also filtered out the adapter sequences (> 10 nt aligned to the adapter, allowing $\leq 10\%$
50 mismatches), low quality reads (i.e. $\geq 10\%$ unidentified nucleotides or $> 50\%$ bases having
51 Phred quality < 5) and duplicated reads generated in the library construction process.

52

53 **Data analysis**

54 ***Reads mapping***

55 The high quality paired-end reads were mapped to the reference chicken genome
56 (Galgal4.78) using Burrows-Wheeler Aligner (BWA) software (version 0.7.8) [1] with the
57 command ‘mem -t 10 -k 32’ and BAM alignment files were generated using SAMtools (version
58 0.1.19) [2].

59 Next, we improved the alignment results by the following steps:

60 (1) The aligned reads with mismatches ≥ 5 or mapping quality = 0 were removed;

61 (2) The alignment results were then corrected using Picard (version 1.96)
62 (<http://broadinstitute.github.io/picard/>) with two core commands. The
63 ‘AddOrReplaceReadGroups’ command was used to replace all read groups in the INPUT file
64 with a new read group and assign all reads to this group in the OUTPUT BAM. The
65 ‘FixMateInformation’ command was used to ensure that all mate-pair information was in sync
66 between each read and its mate pair;

67 (3) Removed potential PCR duplications. If multiple read pairs had identical external
68 coordinates, only the pair with the highest mapping quality was retained;

69 (4) Realigned reads around the InDels. We downloaded variants registered in chicken
70 dbSNP database (Build 145) from NCBI, and generated a target list of intervals by using the
71 command “RealignerTargetCreator” in package Genome Analysis Toolkit (GATK, version 3.1-
72 1- g07a4bf8) [3]. We further used the command “IndelRealigner” to identify regions for
73 realignment where at least one read contains a registered InDel with a cluster of mismatching
74 bases around it.

75 Consequently, ~21.30 Gb high quality data of each individual mapping to reference
76 chicken genome (Additional file 1: Table S1) were used for subsequent analysis.

77

78 ***SNP calling***

79 We first detected individual SNPs simultaneously confirmed by both SAMtools and GATK.
80 The highly accurate alignment was processed using the ‘mpileup’ program in SAMtools with
81 the parameters ‘-C 50 -D -S -m 2 -F 0.002 -d 1000’ (‘-C 50’ is a recommended parameter, ‘-D’
82 and ‘-S’ are default parameters, ‘-m 2’, ‘-F 0.002’ and ‘-d 1000’ are required parameters). The

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

83 variants were then filtered for downstream analysis by requiring a coverage ranging from 4 to
84 200, a minimum root-mean-square mapping quality of 20 and no gaps present within a 3-bp
85 window. Meanwhile, we detected genomic variants for each bird using GATK with the
86 HaplotypeCaller-based method; before calling variants, the base quality scores were
87 recalibrated using command “BaseRecalibrator”, which provides empirically accurate base
88 quality scores for each base in every read. After SNP calling, we applied hard filter command
89 ‘VariantFiltration’ to exclude potential false-positive variant calls with the parameter ‘--
90 filterExpression "QD < 10.0 || FS > 60.0 || MQ < 40.0 || ReadPosRankSum < -8.0" -G_filter
91 "GQ<20"’. As a result, ~6.44 Mb SNPs for each breed/population were identified (Additional
92 file 1: Table S2).

93 Then we merged all individual SNPs into a population SNP-matrix. Finally, we obtained
94 8.53 Mb highly credible SNPs after using strict criteria with filtering MAF (minor allele
95 frequency) < 0.05 and missing genotype > 10% in chicken population. Subsequently, the
96 package ANNOVAR (version May 20, 2013) [4] was used to annotate SNPs causing nonsense
97 and missense mutations.

98 *Insertions and deletions (InDels) calling*

100 The candidate InDels were called along with SNPs by GATK for 91 individuals. We first
101 sifted structural variations for each sample by GATK with the SelectVariants based method.
102 Then, we applied hard filter command ‘VariantFiltration’ to exclude potential false-positive
103 variant calls with the parameter ‘--filterExpression "QD < 2.0 || FS > 200.0 || ReadPosRankSum
104 < -8.0 || InbreedingCoeff < -0.8"’. Finally, we only retained the 1-30 bp InDels for downstream
105 analysis.

106 *Analysis of the population structure and evolutionary history*

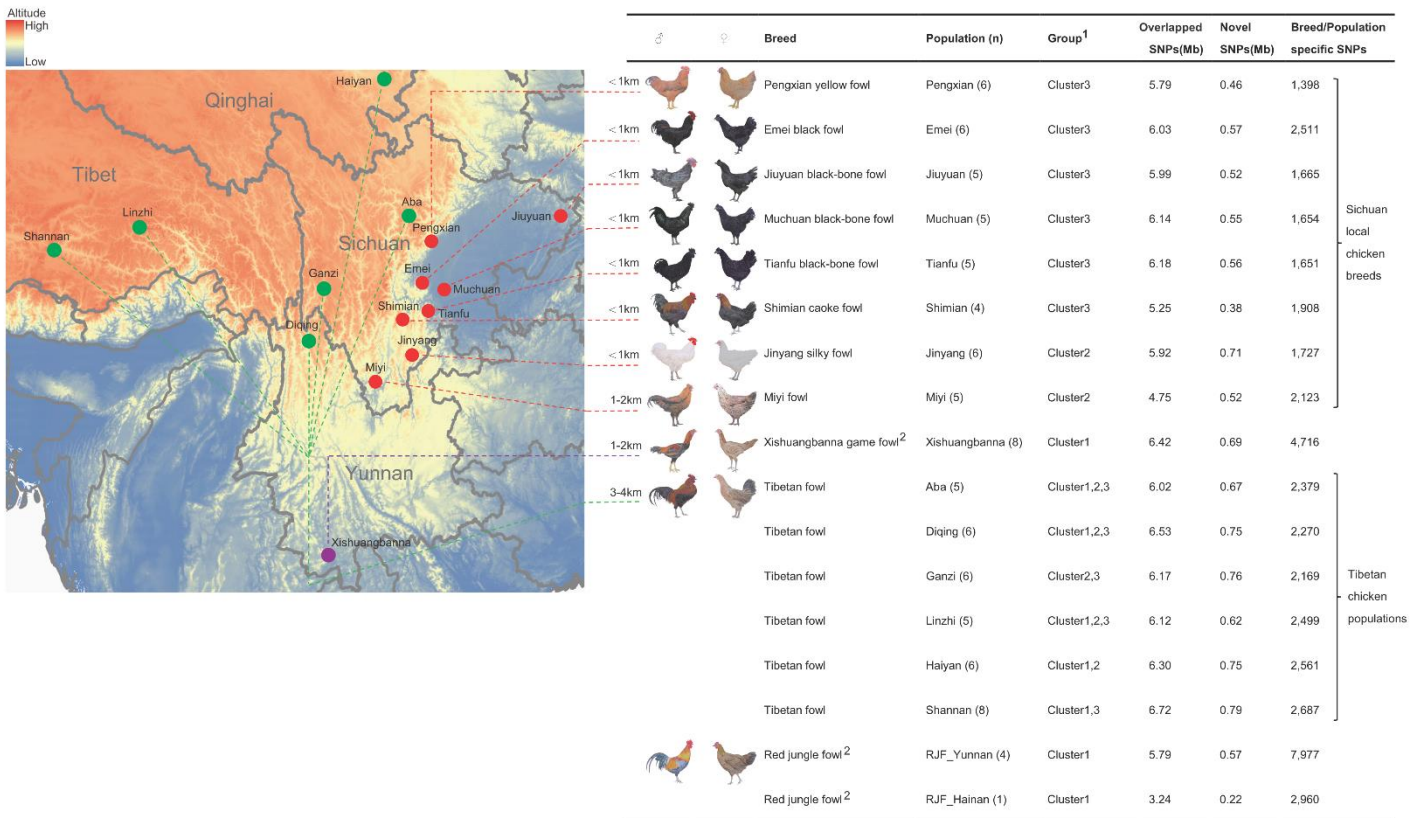
107 Rooted neighbor-joining phylogenetic tree was constructed under the p-distances model in
108 TreeBeST (version 1.9.2) (<http://treesoft.sourceforge.net/treebest.shtml>), using Japanese quail
109 as an outgroup. The reliability of each branch was evaluated by bootstrapping [5] with 1,000
110 replicates. The phylogenetic relationships of the individual genomes were also estimated using
111 principle component analysis (PCA) with the population-scale SNPs using the EIGENSOFT
112 (version 5.0) [6] software, and the eigenvectors were obtained from the covariance matrix
113 generated by R function reigen.

114 **Findings**

115 *Genetic diversity*

116 A total of 7.43 Mb of SNPs out of 8.53 Mb highly credible SNPs were already present in
117 chicken dbSNP database (overlapped SNPs) and 1.10 Mb SNPs were assigned as novel ones.
118 All 1.10 Mb novel SNPs have been submitted to dbSNP (accession numbers from
119 ss2585830405 to ss2586846514 and ss2137077162; see Additional file 2). We further
120
121

122 conducted a comparative genomics analysis of 91 chickens from 15 domestic and 2 wild
 123 populations (Figure 1). The general phenotypic differences between red jungle fowls (RJF),
 124 Tibetan fowls and Sichuan local fowls are shown in Additional file 1: Table S3. We identified
 125 3.46-7.52 Mb SNPs for each breed/population that were confirmed by both SAMtools and
 126 GATK softwares (Additional file 1: Table S2). There were 1,398 to 7,977 SNPs specifically
 127 detected in a breed/population (Figure 1). Nucleotide variability (θ_π) and polymorphism (θ_ω) in
 128 each population were analyzed using the method of sequence diversity statistics [7]. Compared
 129 with Sichuan local chicken breeds ($\theta_\pi = 2.35 \times 10^{-3}$ and $\theta_\omega = 2.13 \times 10^{-3}$), Tibetan chicken
 130 populations have relatively higher genetic diversity ($\theta_\pi = 2.58 \times 10^{-3}$, $P < 2.2 \times 10^{-16}$ and $\theta_\omega =$
 131 2.35×10^{-3} , $P = 0.656$, Mann-Whitney U test) (Additional file 1: Figure S1).



133 **Figure 1.** Sample information and comparison of identified SNPs in each breed/population with
 134 the chicken variants database (dbSNP, Build 145). Overlapped SNPs are SNPs already in
 135 chicken dbSNP. The map displayed here is the geographic distribution of domestic chicken
 136 populations, numbers above the dashed lines are altitudes. Red and green localities represent
 137 eight lowland and six highland chicken populations respectively, sampled in this study. ¹
 138 Individual distribution to each group can be found in Additional file 1: Table S1. ²The whole-
 139 genome sequencing data of eight game fowls and 5 RJFs were downloaded from the NCBI.

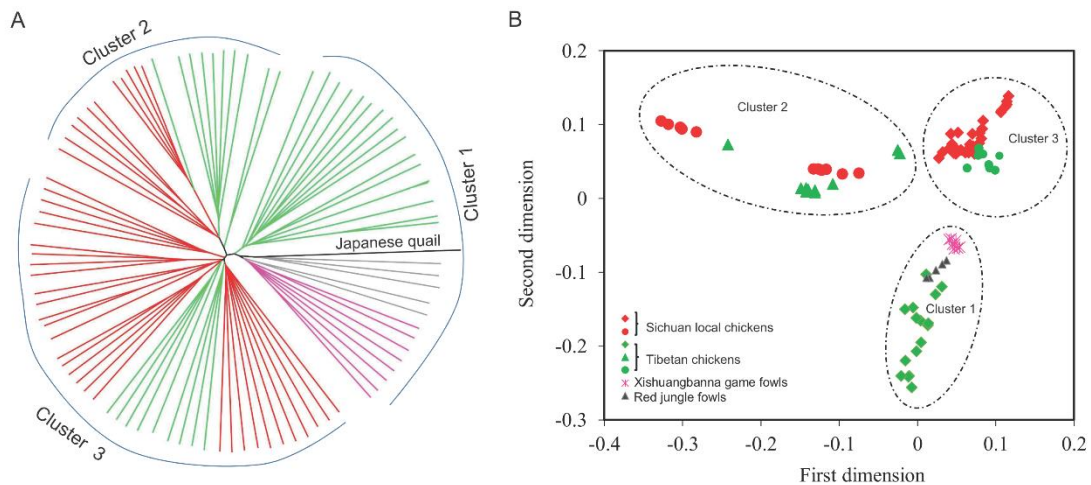
141 As shown in Additional file 1: Figure S2, although most novel SNPs (89.02%) had a low
 142 allele frequency (<0.2 of 91 individuals) compared with the overlapped SNPs (44.02%), only
 143 9,918 (0.88% of 1.10 M) novel SNPs were specifically detected in one breed/population (at
 144 least in an individual). These novel SNPs also exhibited a comparable sequencing depth with
 145 the overlapped SNPs (median of normalized depth of 1.14 versus 1.06) (Additional file 1:

146 Figure S3). In addition, we observed more than 75% of the novel SNPs and 86% of the
147 overlapped SNPs were in non-repeat regions. These results suggest the novel SNPs will serve
148 as a potentially valuable resource for further chicken studies.

149 Overall distribution of the lengths of insertions and deletions (InDels) showed that more
150 than 80% of the InDels were 1-5 bp in length (Additional file 1: Figure S4). Repetitive elements
151 (10.61% of the genome and containing ~15.70% of InDels) are an important source of structural
152 variation in chicken genome (Additional file 1: Figure S5). About a half of InDels (48.39% to
153 51.52%) were occurred in the intergenic regions (588.65 Mb and 56.23% of the genome). The
154 introns (403.35 Mb and containing ~43.86% of InDels) showed higher incidence of InDels than
155 the coding sequences (25.81 Mb and containing ~1.77% of InDels) (Additional file 1: Figure
156 S6). We observed an enrichment of short InDels (1-15 bp in length) in coding sequences that
157 were multiples of 3 bp compared to whole genome sequences, which is expected to preserve
158 the reading frame (Additional file 1: Figure S7).

160 *Population genetics*

161 The neighbor-joining phylogenetic tree revealed the segregation of 15 domestic
162 populations and 2 wild RJF populations into three distinct clusters (cluster 1, cluster 2 and
163 cluster 3) (Figure 2A). A similar pattern of clustering (Figure 2B) was also observed based on
164 principal component analysis (PCA) using EIGENSOFT package [6]. Different from a previous
165 report on the two independent origins of Tibetan chickens [8], we revealed the presence of at
166 least three distinct clusters among the six geographically representative populations of Tibetan
167 fowls: the fowls inhabiting Tibet and Qinghai (in cluster 1) were genetically closer to RJF, while
168 the Tibetan chickens inhabiting Yunnan and Sichuan (cluster 2 and 3) were closer to the
169 domestic populations (Figure 1). These distinct distribution patterns and expansion signatures
170 suggested that the divergent Tibetan clades may have originated from different regions, such as
171 Yunnan, southwest China and/or surrounding areas [8]. We found that many Tibetan chickens
172 clustered with other Sichuan local chicken breeds in cluster 2 and cluster 3, which may be
173 attributable to shared ancestral polymorphism and/or recent introgression events by way of
174 possible crossbreeding between Tibetan chicken with the geographically neighboring Sichuan
175 local chickens. Although this inference is consistent with recent breeding activities in Tibet
176 plateau [8], further analysis are required to explore the introgression between them.



177
178 **Figure 2. Population genetics of studied chickens.** (A) Rooted neighbor-joining phylogenetic

179 tree with the neighbor-joining method, using Japanese quail as an outgroup. The reliability of
180 each branch was evaluated by bootstrapping with 1,000 replicates. Different groups of chicken
181 populations: Sichuan local chickens (red), Tibetan chickens (green), the Xishuangbanna game
182 fowls (purple), RJFs (grey) and Japanese quail (black). (B) Principal component plots. The first
183 dimension and second dimension are shown. The fraction of the variance explained was 8.91%
184 for eigenvector 1 ($P<0.05$, Tracy-Widom test) and 7.43% for eigenvector 2 ($P<0.05$, Tracy-
185 Widom test).

186

187 **Conclusion**

188 Understanding the nature of diversifying selection, especially detecting selection
189 signatures, and identifying genes in a genome that are, or have been, under selection have been
190 the hot topics of interests. This study provides comparative genomic landscape of variations in
191 17 chicken populations to understand genetic variations underlying the phenotypic diversity of
192 chicken breeds/populations. This data will serve as a valuable resource for investigating
193 diversifying selection and candidate genes for selective breeding in chicken.

194

195 **Availability of supporting data**

196 The sequencing data for this project have been deposited in the NCBI sequence read archive
197 (SRA) under accession number SRP067615. All supplementary Figures and Tables are
198 provided in Additional file 1.

199

200 **Additional file**

201 Additional file 1: Table S1, Table S2, Table S3, Figure S1, Figure S2, Figure S3, Figure S4 and
202 Figure S5. (doc 1.3 MB). **Table S1.** A summary of the chickens used in this study: regions of
203 collection/popularization and coverage and mean depth of resequencing. **Table S2.** SNPs
204 annotation and genetic diversity of 17 chicken populations analyzed in this study. **Table S3.**
205 The general phenotypic differences between red jungle fowls, Tibetan and Sichuan local
206 chickens. **Figure S1.** Average nucleotide polymorphism (θ_w) and nucleotide diversity (θ_π)
207 among Sichuan local chickens, Tibetan chickens and red jungle fowls. **Figure S2.** Allele
208 frequency spectra in 91 birds and Number of alleles distribute in 1 to 17 chicken
209 breeds/populations. **Figure S3.** Comparison of sequencing depth between the SNPs that are
210 already in dbSNP (overlapped SNPs) and novel SNPs. **Figure S4.** Overall distribution of the
211 lengths of InDels (1-30 bp). **Figure S5.** Percentage composition of InDels in repeat elements.
212 **Figure S6.** Percentage distribution (A) and probability (B) for InDels across different genomic
213 elements. **Figure S7.** Length distribution of small InDels in the whole genome (A) and coding
214 sequence (CDS) regions (B).

215 Additional file 2: Accession numbers of 1.10 Mb novel SNPs. (txt 20.3 Mb)

216

217 **Funding**

218 This work was supported by China Agricultural Research System (CARS-41), the 12th Five
219 Year Plan for breeding program in Sichuan-Selective breeding of new breeds and the synthetic
220 strains in laying hens (2011NZ0099-7), National Natural Science Foundation of China

221 (31402063), Sichuan Provincial Department of Science and Technology Program
222 (2015JQ0023) and the National Program for Support of Top-notch Young Professionals and
223 the Young Scholars of the Yangtze River.

224 **Authors' contributions**

225 Q.Z., and MZ.L. designed and supervised the project. B.C., M.L., H.Y., Y.W., X.Z., G.Z., U.G.,
226 MJ.L., L.Z., M.Y., R.J., R.L., and X.Z collected and generated the data, and performed the
227 preliminary bioinformatic analyses. T.C., S.T., Y.L., Z.X., L.J., Q.T., H.X., and X.Z. filtered the
228 data and performed the majority of the population genetic analysis. D.L. and T.C. wrote the
229 manuscript.

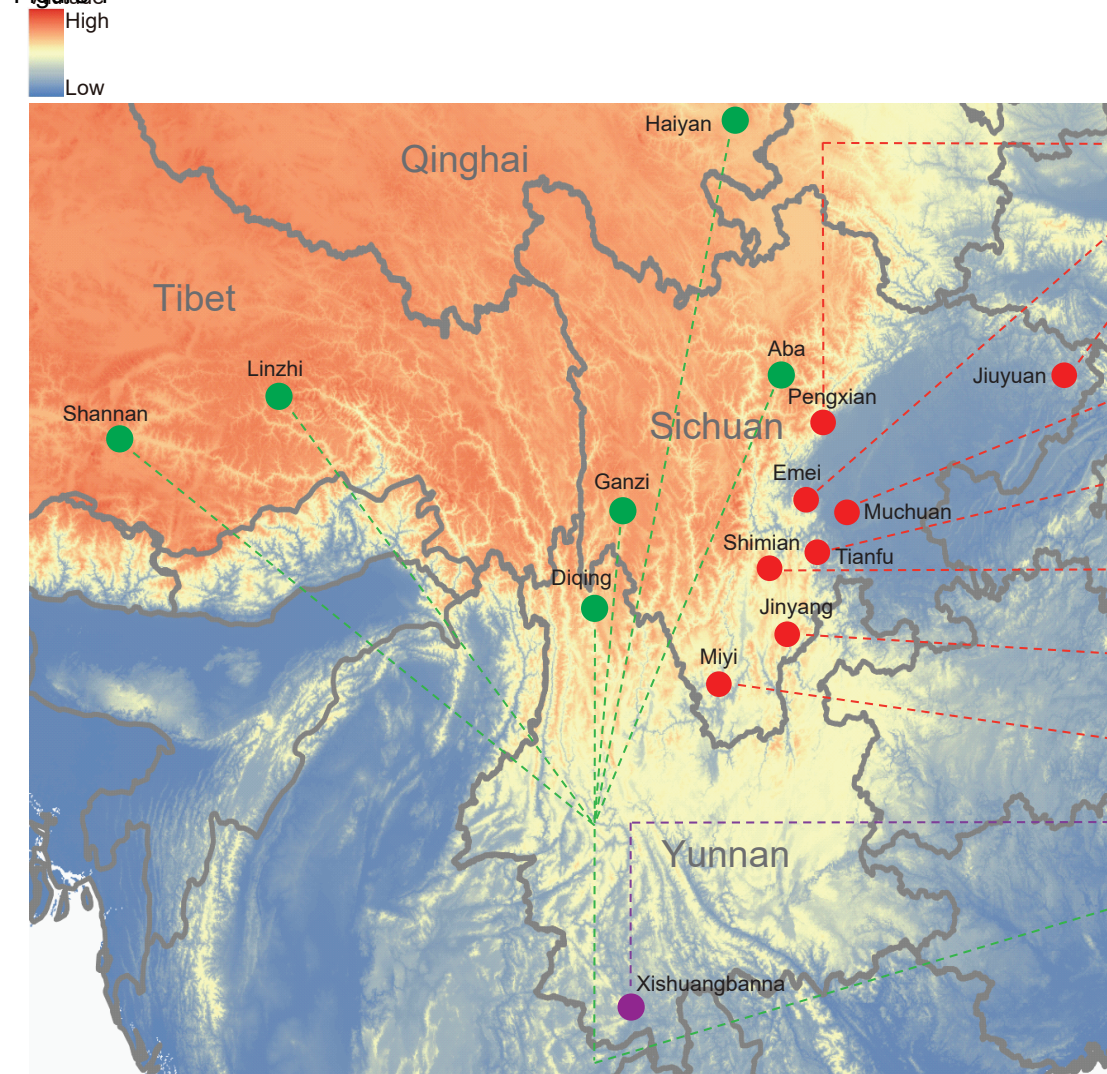
231 **Competing financial interests**

232 The authors declare no competing financial interests.

234 **References**

235
236 1. Li H, Durbin R. Fast and accurate long-read alignment with Burrows–Wheeler
237 transform. *Bioinformatics*. 2010; 26: 589-595.
238 2. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The sequence
239 alignment/map format and SAMtools. *Bioinformatics*. 2009; 25: 2078-2079.
240 3. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, et al. The
241 Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA
242 sequencing data. *Genome Res*. 2010; 20: 1297-1303.
243 4. Conesa A, Götz S, García-Gómez JM, Terol J, Talón M, Robles M. Blast2GO: a universal
244 tool for annotation, visualization and analysis in functional genomics research.
245 *Bioinformatics*. 2005; 21: 3674-3676.
246 5. Felsenstein J. Confidence limits on phylogenies: an approach using the bootstrap.
247 *Evolution*. 1985; 39: 783-791.
248 6. Price AL, Patterson C, Plenge RM, Weinblatt M, Shadick NA, Reich D. Population structure and
249 eigenanalysis. *PLoS Genetics*. 2006; 2: 2074--2093.
250 7. Nei M, Li WH. Mathematical model for studying genetic variation in terms of
251 restriction endonucleases. *Proc Natl Acad Sci U S A*. 1979; 76: 5269-5273.
252 8. Ming-Shan W, Yan L, Min-Sheng P, Li Z, Zong-Ji W, Qi-Ye L, et al. Genomic Analyses
253 Reveal Potential Independent Adaptation to High Altitude in Tibetan Chickens.
254 *Molecular Biology & Evolution*. 2015; 32: 1880-9.

Figure 1

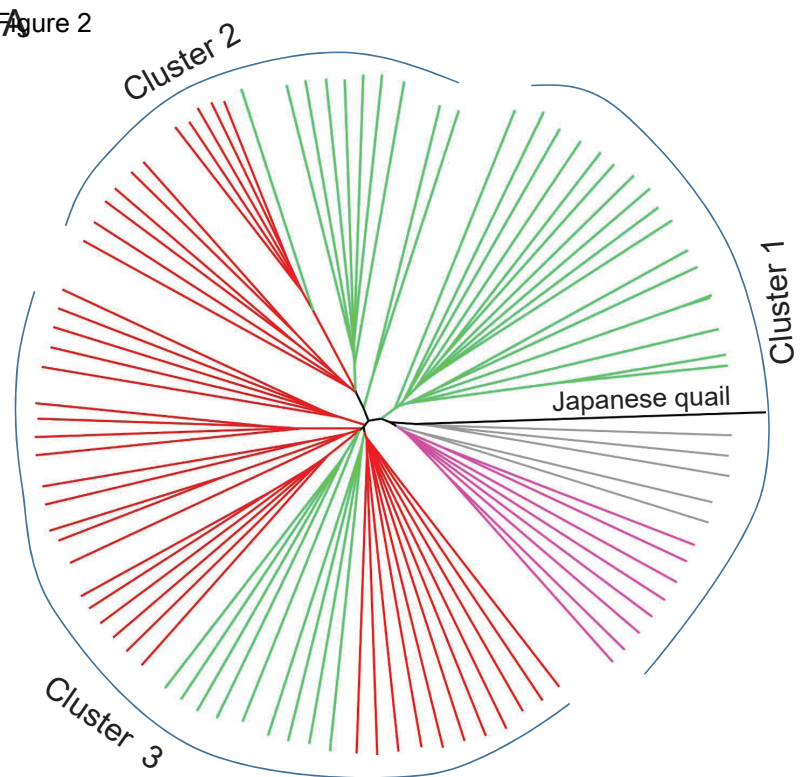
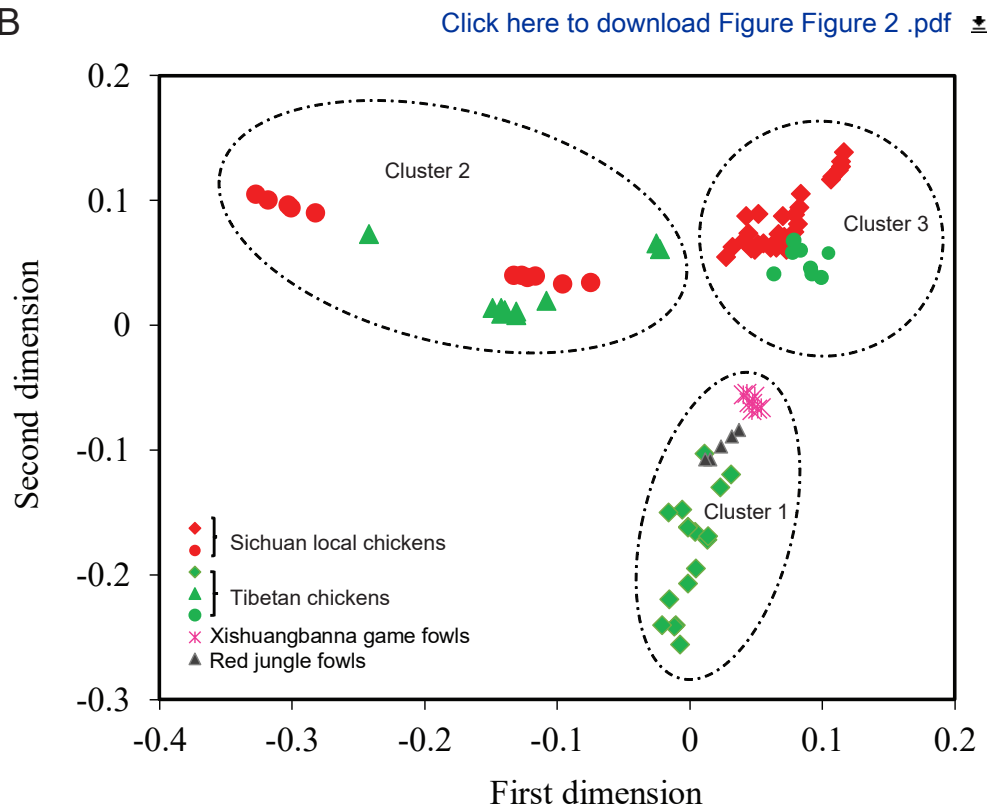
[Click here to download Figure Figure 1.pdf](#)


♂	♀	Breed	Population (n)	Group ¹	Overlapped SNPs(Mb)	Novel SNPs(Mb)	Breed/Population specific SNPs
		Pengxian yellow fowl	Pengxian (6)	Cluster3	5.79	0.46	1,398
		Emei black fowl	Emei (6)	Cluster3	6.03	0.57	2,511
		Jiu yuan black-bone fowl	Jiu yuan (5)	Cluster3	5.99	0.52	1,665
		Muchuan black-bone fowl	Muchuan (5)	Cluster3	6.14	0.55	1,654
		Tianfu black-bone fowl	Tianfu (5)	Cluster3	6.18	0.56	1,651
		Shimian caoke fowl	Shimian (4)	Cluster3	5.25	0.38	1,908
		Jinyang silky fowl	Jinyang (6)	Cluster2	5.92	0.71	1,727
		Miyi fowl	Miyi (5)	Cluster2	4.75	0.52	2,123
		Xishuangbanna game fowl ²	Xishuangbanna (8)	Cluster1	6.42	0.69	4,716
		Tibetan fowl	Aba (5)	Cluster1,2,3	6.02	0.67	2,379
		Tibetan fowl	Diqing (6)	Cluster1,2,3	6.53	0.75	2,270
		Tibetan fowl	Ganzi (6)	Cluster2,3	6.17	0.76	2,169
		Tibetan fowl	Linzhi (5)	Cluster1,2,3	6.12	0.62	2,499
		Tibetan fowl	Haiyan (6)	Cluster1,2	6.30	0.75	2,561
		Tibetan fowl	Shannan (8)	Cluster1,3	6.72	0.79	2,687
		Red jungle fowl ²	RJF_Yunnan (4)	Cluster1	5.79	0.57	7,977
		Red jungle fowl ²	RJF_Hainan (1)	Cluster1	3.24	0.22	2,960

Sichuan local chicken breeds


Tibetan chicken populations

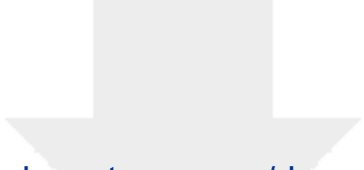
Figure 2

**B**



Click here to access/download
Supplementary Material
Additional file 1-R2.docx





Click here to access/download
Supplementary Material
Additional file 2.txt

