

1 The chicken pan-genome reveals gene content variation and a
2 promoter region deletion in *IGF2BP1* affecting body size

3 Kejun Wang^{1,2†}, Haifei Hu^{3†}, Yadong Tian^{1,2}, Jingyi Li⁴, Armin Scheben⁵, Chenxi
4 Zhang^{1,2}, Yiyi Li^{1,2}, Junfeng Wu^{1,2}, Lan Yang^{1,2}, Xuwei Fan^{1,2}, Guirong Sun^{1,2},
5 Donghua Li^{1,2}, Yanhua Zhang^{1,2}, Ruili Han^{1,2}, Ruirui Jiang^{1,2}, Hetian Huang^{1,2}, Fengbin
6 Yan², Yanbin Wang², Zhuanjian Li^{1,2}, Guoxi Li^{1,2}, Xiaojun Liu^{1,2}, Wenting Li^{1,2*}, David
7 Edwards^{3*}, Xiangtao Kang^{1,2*}

8 ¹College of Animal Science and Technology, Henan Agricultural University,
9 Zhengzhou 450046, China

10 ²Henan Key laboratory for innovation and utilization of chicken germplasm
11 resources, Zhengzhou, 450046, China

12 ³School of Biological Sciences and Institute of Agriculture, University of Western
13 Australia, Crawley, 6009 WA, Australia

14 ⁴Key Laboratory of Agricultural Animal Genetics, Breeding and Reproduction of
15 Ministry of Education, College of Animal Science and Technology, Huazhong
16 Agricultural University, 430070 Wuhan, Hubei, China

17 ⁵Simons Center for Quantitative Biology, Cold Spring Harbor Laboratory, Cold Spring
18 Harbor, NY, USA

19 †These authors contributed equally to this work.

20 *Correspondence and requests for materials should be addressed to W. L.
21 (email:liwenting_5959@hotmail.com), D.E. (email: dave.edwards@uwa.edu.au) and
22 X.K. (email:xtkang2001@263.net).

23 **Abstract**

24 Domestication and breeding have reshaped the genomic architecture of chicken, but the
25 retention and loss of genomic elements during these evolutionary processes remain
26 unclear. We present the first chicken pan-genome constructed using 664 individuals,
27 which identified an additional ~66.5 Mb sequences that are absent from the reference
28 genome (GRCg6a). The constructed pan-genome encoded 20,491 predicated protein-
29 coding genes, of which higher expression level are observed in conserved genes relative
30 to dispensable genes. Presence/absence variation (PAV) analyses demonstrated that
31 gene PAV in chicken was shaped by selection, genetic drift, and hybridization. PAV-
32 based GWAS identified numerous candidate mutations related to growth, carcass
33 composition, meat quality, or physiological traits. Among them, a deletion in the
34 promoter region of *IGF2BP1* affecting chicken body size is reported, which is
35 supported by functional studies and extra samples. This is the first time to report the
36 causal variant of chicken body size QTL located at chromosome 27 which was
37 repeatedly reported. Therefore, the chicken pan-genome is a useful resource for
38 biological discovery and breeding. It improves our understanding of chicken genome
39 diversity and provides materials to unveil the evolution history of chicken
40 domestication.

41

42 **Introduction**

43 Chicken (*Gallus gallus*) is the most abundant domesticated animal in the world. The
44 publication of the chicken genome in 2004 (Hillier, et al. 2004) paved the way to
45 identify the QTLs or QTGs involved in economically important traits, dissect the
46 evolutionary processes of domestication, and understand the genetic basis of distinct
47 phenotypes differentiating domesticated chickens and their wild relatives. Recently, the
48 domestic chicken *Gallus gallus domesticus* was reported to have been domesticated
49 from one subspecies of red jungle fowl, *Gallus gallus spadiceus* (Wang, Thakur, et al.
50 2020b). Nevertheless, subspecies of *Gallus gallus* and other jungle fowls can introgress
51 with *Gallus gallus domesticus* and these interspecies hybridizations have affected the

52 genetic content of the species during evolution (Barton 2001; DESTA 2019; Lawal, et
53 al. 2020; Wang, Thakur, et al. 2020b). Traits such as yellow skin, pencilled feathers and
54 the spotted comb of domesticated chickens are likely the result of introgressions from
55 *Gallus sonneratii*, *Gallus varius* and *Gallus lafayettei* (Morejohn 1968b; Eriksson, et al.
56 2008; Fallahshahroudi, et al. 2019). Hybridizations leading to fertile offspring have
57 been documented between *Gallus* species (DANFORTH 1958; Morejohn 1968a).
58 These indicate that *Gallus gallus domesticus* is an admixed species, not only derived
59 from red jungle fowl (Wang, Thakur, et al. 2020a). A recent study also found different
60 genome sizes between red jungle fowl and domestic chicken lineages (Piegu, et al.
61 2020). Moreover, growing evidence suggests that structural variations are present in a
62 substantial proportion of the genomes of many animals (Bickhart and Liu 2014),
63 including human (Sherman, et al. 2019), pig (Zhao, et al. 2016; Li, Chen, et al. 2017;
64 Tian, et al. 2020), salmon (Bertolotti, et al. 2020), and chicken (Kerstens, et al. 2011;
65 Seol, et al. 2019). A range of phenotypes in chicken was reported to be determined by
66 structural variations, such as feathered legs (Li, Lee, et al. 2020), crest (Li, et al. 2021),
67 blue egg shell (Wang, et al. 2013), muffs and beard (Guo, et al. 2016), comb (Wright,
68 et al. 2009; Imsland, et al. 2012), and fibromelanosis (Dorshorst, et al. 2011). The
69 current chicken reference genome (GRCg6a) is derived from a single red jungle fowl
70 individual. This reference therefore cannot fully capture the genetic diversity of
71 domesticated chickens, and may be unable to reveal the genetic basis of some
72 phenotypes. Recently, an increasing number of reports for pan-genomes in human
73 (Sherman, et al. 2019), pig (Tian, et al. 2020), goat (Li, et al. 2019), and also in plants
74 (Bayer, et al. 2020), have focused on capturing genetic variations between different
75 individuals within the species. The pan-genome represents the gene set of the species
76 rather than a representative individual, which can uncover the genetic diversity and
77 resolve structural variations that are missed by studies using a single reference genome.
78 The pan-genome can also provide a straightforward way to detect presence/absence
79 variations (PAV) and explore the distributions of these variants at the population level.

80 Body size is an important quantitative trait that has been intensively selected during

81 chicken improvement and possibly associated with genome structural variations. One
82 of the well-known candidate genes linked to body size is insulin-like growth factor 2
83 mRNA-binding protein 1 (*IGF2BP1*). *IGF2BP1* can regulate cell proliferation,
84 differentiation, morphology, and metabolism through regulating mRNA localization,
85 stability, and translation of targeted genes (Stohr, et al. 2012; Bell, et al. 2013). In recent
86 studies, *IGF2BP1* was reported as N6-methyladenosine (m6A) readers to regulate the
87 above functions (Huang, et al. 2018; Zhu, et al. 2020; Zhang, et al. 2021). Knockout of
88 *IGF2BP1* in mouse led to mild active colitis, mild-to-moderate active enteritis, and
89 decreasing of barrier function and body weight (Singh, et al. 2020). Dwarfism and
90 impaired gut development were also observed in *IGF2BP1*-deficient mice (Hansen, et
91 al. 2004). Evidence from genome-wide association studies (GWAS) and QTL mapping
92 revealed that the genomic regions upstream of *IGF2BP1* were significantly associated
93 with body weight, head weight, gizzard weight, chest width, leg weight, and wing
94 weight in chicken and duck (Sheng, et al. 2013; Ma, et al. 2018; Zhou, et al. 2018;
95 Wang, Bu, et al. 2020; Wang, Cao, et al. 2020; Zhang, et al. 2020). However, the causal
96 variations of *IGF2BP1* that responsible for body sizes in chicken and duck remain
97 unclear.

98 Here, we constructed the first chicken pan-genome and comprehensively
99 investigate PAV using this pan-genome, revealing changes in allele frequencies
100 associated with chicken evolution. We found that deletions in the promoter region of
101 *IGF2BP1* can increase transcriptional activity and gene expression, regulating the body
102 size in commercial chickens. Dissection of the causal variation of *IGF2BP1* associated
103 with body size can accelerate the breeding process for high growth rate chickens using
104 marker-assisted selection. These findings will improve our understanding of changes in
105 chicken gene content during domestication and breeding and help to design highly
106 productive chicken breeds in the future

107 **Results**

108 **Pan-genome construction of chicken**

109 We constructed the first *Gallus gallus* pan-genome using an iterative mapping and

110 assembly approach based on the chicken reference genome GRCg6a assembly. A set of
111 whole-genome sequencing data including 664 individuals was used in the pan-genome
112 construction, which contains 5 *Gallus gallus* wild subspecies, 28 native breeds
113 (indigenous chicken breeds raised by farmers that did not experience intense artificial
114 selection) and 4 commercial breeds (Supplementary Table S1; Figure 1a).

115 The *Gallus gallus* pan-genome identified an additional ~66.5 Mb sequences that
116 are absent from the reference genome (GRCg6a), encoding an additional 4,063 high-
117 confidence genes (Supplementary Table S2-S3). Of these, 49% (1,976 genes) non-
118 reference genes are only present in a small proportion of chickens (Figure 1b). Together,
119 the chicken pan-genome, including reference and non-reference sequences, consists of
120 1131.9 Mb and contains 20,941 predicted protein-coding genes. A total of 81 RNA-seq
121 datasets from 27 tissues (including digestive, respiratory, kinetic, urinary, reproductive,
122 endocrine, circulatory, nervous, immune, epithelium, and connective system) were used
123 to investigate the gene expression (Supplementary Table S4). We observed an average
124 normalized transcript per million (TPM) abundance greater than 1 for 90.6% of the
125 autosomal genes in the reference genome and 19.4% of the non-reference genes. This
126 pattern is similar to those found in other pan-genome studies (Zhao, et al. 2018; Gao,
127 et al. 2019), which showed that genes in the reference genome generally have higher
128 expression than genes in the non-reference contigs (Supplementary Figure S1a).

129 **Discovery of gene Presence/Absence Variation (PAV)**

130 After sample selection (see Methods and Supplementary Figure S2), a total of 268
131 individuals with average sequence depth larger than 10x based on pan-genome
132 estimation were available for gene PAV detection, including 6 wild, 217 native and 45
133 commercial individuals (Supplementary Table S1).

134 We categorized genes in the chicken pan-genome according to their gene presence
135 frequencies. 15,205 (76.32%) core genes are shared by 268 individuals. 4,738 genes are
136 variable including 391 softcore, 2,351 shell and 1,976 cloud genes, which are present
137 in more than 99%, 1-99%, and less than 1% of all individuals, respectively (Figure 1b).
138 The chicken pan-genome showed a moderate core gene content (76.32%) compared to

139 that of human (96.88%) (Duan, et al. 2019), mussel (69.2%) (Gerdol, et al. 2020) , and
140 plants (35 ~81%) (Gao, et al. 2019). Gene Ontology (GO) enrichment results of each
141 cluster of variable genes are presented in Supplementary Table S5. Variable genes were
142 enriched in the function associated with reproduction, nutrient absorption, metabolic
143 and biosynthetic process (Figure 1c). RNA-seq analysis revealed that the expression
144 level of flexible genes (shell and cloud genes) was significantly lower than that of
145 conserved genes (core and softcore genes) (Supplementary Figure S1b). No apparent
146 difference of expression was identified between conserved genes in the reference and
147 non-reference sequences, but the expression of conserved genes was significantly
148 higher than that of flexible genes in both reference and non-reference sequences
149 (Supplementary Figure S1c). Pan-genome modelling revealed a closed pan-genome
150 with an estimated total of 19,190 genes (genes on sex chromosomes were excluded
151 because the gene content was different between chromosomes Z and W) (Figure 1d).
152 This suggests the chicken pan-genome assembled using our selected 268 individuals
153 included all or nearly all of the *Gallus gallus* gene content.

154 **Gene PAV shaped by selection, genetic drift and hybridization**

155 We observed a broad gene PAV distribution within different groups, with substantial
156 variation in the native chickens and wild relatives (red jungle fowls) (Figure 2a). PAV-
157 based PCA and phylogenetic analysis also showed high diversity among wild relatives
158 and native chickens, while commercial broiler and layer clustered together (Figure 2b-
159 c). Moreover, two clades of commercial chickens were further separated into two
160 groups, meat-production (two broiler breeds: BRA and BRB) and egg-production (two
161 layer breeds: BL and WL). These differences between commercial and native or wild
162 chickens are likely due to selection, but the genetic drift and other factors can not be
163 ruled out. Therefore, we further investigated whether selection, genetic drift, and
164 hybridization can alter gene PAV content since SNP-based allelic frequencies can be
165 shaped by selection, genetic drift, and hybridization (Edwards 2008).

166 We analyzed the pool sequencing data of ‘Virginia body weight lines’ and
167 compared the gene PAV content between high weight selected (HWS) lines and low

168 weight selected (LWS) lines which were divergently selected from the same founder
169 White Plymouth Rock population (Lillie, et al. 2018). Two lines had been suffered from
170 intensive bidirectional selection for 8-week body weight, between which about 15-fold
171 phenotypic difference presented. PAV-based PCA and phylogenetic analysis showed
172 two distinct clusters were consistent with selected lines (Supplementary Figure S3a-3b).
173 This suggests that gene PAV can be shaped by intensively artificial selection. We have
174 further compared the frequencies of gene PAV between HWS and LWS and identified
175 the candidate genes related to the intensive bidirectional selection for body weight.
176 Twenty-four genes were found to be completely absent in HWS and present in LWS or
177 entirely present in HWS and absent in LWS (Supplementary Figure S3c). The well
178 studied *SH3 domain containing ring finger 2* (*SH3RF2*, *ENSGAT00000090177*) gene,
179 regulating appetite and affecting body weight, was also identified as one of these genes
180 that is completely absent in HWS but present in LWS (Rubin, et al. 2010; Jing, et al.
181 2020).

182 We further investigated whether gene PAV within the chicken population can be
183 affected by genetic drift or hybridization. Firstly, we studied conserved populations of
184 varying size. The subpopulations GS1, GS2, and GS3 are from the Gushi chicken
185 populations, of which GS1 (n=30) and GS2 (n=30) were sampled from a small
186 conserved population in 2010 and 2019 respectively, while GS3 (n=30) was sampled
187 from a large conserved population in 2019. The subpopulations XB1 (n=30) and XB2
188 (n=30), are the Xichuan black-bone chicken populations, which were sampled from a
189 large conserved population in 2010 and 2019, respectively (Supplementary Note;
190 Supplementary Figure S4a). We did not observe the change of PAV content during short
191 period (less than 9 years), whatever in a small or big conservation population, by
192 comparing XB1 and XB2, or GS1 and GS2. However, we observed an apparent division
193 between GS1 or GS2 and GS3 based on the results of PCA and phylogenetic analysis
194 (Supplementary Figure S4b-d). We also found a significant reduction of genetic
195 diversity in GS1 and GS2 in comparison with GS3 based on SNP heterozygosity and
196 allelic richness analysis (observed heterozygosity: 0.18 in GS1, 0.17 in GS2, 0.23 in

197 GS3; allelic richness: 0.53 in GS1, 0.47 in GS2, 0.75 in GS3). These results are
198 consistent with significant differences in gene PAV content (Supplementary Figure S4b-
199 d) and previous studies showing that small conserved populations suffer from genetic
200 drift after long periods of isolation which leads to a reduction of genetic diversity
201 (Whitlock 2000). Based on the above evidence, we compared the gene PAV frequencies
202 between GS3 and GS1+GS2 to investigate gene PAV involving genetic drift by long
203 periods of isolation. According to Fisher's exact test (FDR < 0.001) (Gao, et al. 2019),
204 we only identified six genes that were significantly different in frequencies between
205 GS3 and GS1+GS2 (Supplementary Figure S4e), and none of these genes has a clear
206 functional annotation (annotated as proteins without known function). Of these, four
207 gene PAVs were fixed or nearly fixed in GS1+GS2, that consistent with the reduction
208 of genetic diversity of GS1 and GS2. Secondly, we compared the gene PAV between
209 Gushi chickens and the Gushi×Anak F2 population (Supplementary Figure S5a) and
210 identified a clear divergence between Gushi breeds and the F2 population according to
211 the PAV distribution (Supplementary Figure S5b). PAV-based PCA and the
212 phylogenetic analysis also revealed that Gushi pure breeds and hybrid population fall
213 into two distinct clades. These results suggest a relatively larger effect of hybridization
214 on gene content, which is also significantly more extensive than that from genetic drift
215 by comparing their clustering (Supplementary Figure S5c-d).

216 **Change of gene PAV frequency during breeding**

217 Gene PAV can be shaped by domestication and improvement; therefore, PAV within
218 populations can also be applied to track the evolutionary history of a species (Gao, et
219 al. 2019; Guo, et al. 2020). By comparing the gene presence frequency between the
220 commercial and native breeds, we identified 30 significantly increased genes and 83
221 significantly decreased genes associated with post-domestication breeding
222 improvement (Supplementary Table S6; Supplementary Figure S6). Of these, 10
223 significant genes (7 increased and 3 decreased) are located on the reference genome.
224 We observed that two uncharacterized genes (*PanGallus_Gene02610* and
225 *ENSGALT00000098327*) are lost in modern breeds. We also observed four immune-

226 related genes significantly decreased during improvement, including a class I
227 histocompatibility antigen (*ENSGALT00000081489*), a B-cell differentiation antigen
228 CD72-like (*PanGallus_Gene00218*), a T-cell differentiation antigen CD6
229 (*PanGallus_Gene04583*), and an Immunoglobulin G-binding protein A
230 (*PanGallus_Gene03891*).

231 Tibetan chicken living at the Tibetan Plateau shows the environmental adaptation
232 to high altitudes, particularly to the hypoxic environment (Wang, et al. 2015). Therefore,
233 we compared the gene PAV frequencies between Tibetan chicken and other lower land
234 indigenous chicken to identify candidate genes associated with the environmental
235 adaptation to high altitude (Supplementary Table S7). A total of 121 genes showing
236 significant difference in PAV frequencies were identified, of which frequencies of 118
237 genes were significantly increased in Tibetan chicken. *Vasodilator-stimulated*
238 *phosphoprotein* (*VASP*, *ENSGALT00000100137*) was found to have a high presence
239 frequency (0.906) in Tibetan chicken compared to other lower land native chickens
240 (0.476). *VASP* has been reported to protect endothelial barrier function during hypoxia
241 (Schmit, et al. 2012). Vasculature of *VASP* deletion mouse exhibited patterning defects
242 and lacks structural integrity, leading to edema and hemorrhaging (Furman, et al. 2007).
243 This evidence suggests *VASP* is likely to play an essential role in vasculature function
244 and structure in a hypoxic environment. *Transitional endoplasmic reticulum ATPase*
245 *gene* (*ENSGALT00000056168*) was nearly completely lost in Tibetan chicken
246 (frequency is 0.093), while had moderate frequency in other lower land native chickens
247 (0.568). Previous studies revealed that transitional endoplasmic reticulum ATPase
248 activity is significantly inhibited during hypoxia in rat and western painted turtles
249 (Henrich and Buckler 2013; Smith, et al. 2015). This suggests that the absence of
250 *transitional endoplasmic reticulum ATPase* gene is potentially associated with the
251 adaptation to hypoxic environment.

252 **Change of PAV frequency in promoter regions during breeding**

253 Most PAV analysis in previous pan-genome studies has focused on the protein coding
254 regions. However, further investigations of the roles of regulatory regions are also

255 required since they can affect gene expression and phenotype (Van Laere, et al. 2003;
256 Swinnen, et al. 2019). Similarity between orthologous promoters drastically decreased
257 when distance was longer than 2 kb from the gene transcription start site (TSS)
258 (Keightley, et al. 2005). Therefore, promoter regions are generally anchored within the
259 2 kb upstream genomic region of the TSS (Farre, et al. 2007; Abe and Gemmell 2014).
260 In this study, the promoter region was defined as the 3 kb upstream genomic region
261 from TSS to maximize the captured promoter regions. In order to detect smaller PAV
262 in the promoter region, we divided each of the promoter region into three windows: 0-
263 1 kb, 1-2 kb, and 2-3 kb upstream of the gene and investigated the frequencies of PAV
264 in each window (Figure 3). We observed that frequencies of 143 PAVs in the 0-1 kb
265 region of commercial chickens were significantly different from that of native chickens,
266 which contains 117 increased and 26 decreased. In the same comparison, the
267 frequencies of 80 PAVs differed significantly in the 1-2 kb regions with 56 increased
268 and 24 decreased, and 78 PAVs differed in the 2-3 kb regions with 55 increased and 23
269 decreased (Figure 3a-c; Supplementary Table S8). We found 12 genes in the olfactory
270 receptor gene family that showed reduced presence frequency in the promoter regions
271 of commercial chickens relative to native chickens (Supplementary Figure S7a). We
272 also observed that the presence frequencies of the promoter region of 9 immunoglobulin
273 related genes were altered during improvement (Supplementary Figure S7b). Genes
274 with significantly altered PAVs frequencies in promoter regions during breeding were
275 enriched mainly in the GO terms of modulation by virus of host process, cyclin-
276 dependent protein kinase holoenzyme complex, and p53 binding (Supplementary
277 Figure S8).

278 Interestingly, we found two PAVs located at both 1-2 kb and 2-3 kb upstream region
279 of *IGF2BP1* gene respectively, which their presence frequencies are significantly less
280 in commercial chickens than native chickens (Figure 3b-c). A high loss rate was
281 observed in commercial breeds compared to native breeds, with a 1-2 kb promoter
282 region presence frequency of 0.04 in commercial breeds and 0.83 in native breeds
283 (Supplementary Table S8). Similarly, the 2-3 kb promoter region presence frequency

284 was 0.04 in commercial breeds and 0.89 in native breeds (Figure 3b-c).

285 **PAV-based GWAS on promoter regions**

286 To uncover traits determined by promoter region PAV, we further conducted PAV-based
287 GWAS on the promoter regions using the Gushi×Anak population with 204 F2
288 individuals (Figure 3d-f). Anak chicken is a commercial broiler breed from Israel, while
289 Gushi is an indigenous chicken of China that did not experience from an intensive
290 selection. We identified 56 association events for 0-1 kb promoter regions, 61 for 1-2
291 kb promoter regions and 78 for 2-3 kb promoter regions (Supplementary Table S8).
292 These association events are involved in 81 traits, including body size, growth, carcass,
293 meat quality, and physiological traits (Supplementary Note). For example, the PAV for
294 2-3 kb upstream region of *ENSGALG00000052768* (low-density lipoprotein receptor
295 precursor, *LDLR*) was functionally related to serum CREA (creatinine) level.
296 *ENSGALG00000051173* (olfactory receptor 14C36-like) was found to be associated
297 with ileum length (IL), jejunum length (JL) and cecum length (CL). We also found that
298 the promoter region PAV of immune-related genes showed associations with production
299 traits. For instance, *ENSGALG00000054397* (class I histocompatibility antigen, F10
300 alpha chain-like isoform X1) was associated with breast bone length (BBL12) and
301 *ENSGALG00000050329* (class I histocompatibility antigen, F10 alpha chain-like
302 isoform X1) was correlated with body weight at birth (BBW). *ENSGALG00000051088*
303 (Gallus gallus class I histocompatibility antigen, F10 alpha chain-like) was linked with
304 BBL12 and body slanting length at 12 weeks (BSL12) (Supplementary Table S8).
305 Immunoglobulin related genes were also identified to correlate with production traits.
306 *ENSGALG00000049846* (immunoglobulin-like receptor CHIR2D-751 precursor) was
307 associated with breast muscle weight (BMW) and the ratio of head weight to body
308 weight at 12 weeks (HR1). *ENSGALG00000045164* (leukocyte immunoglobulin-like
309 receptor subfamily A member 2) was associated with breast muscle weight (BMW) and
310 shank girth (SG8). *ENSGALG00000050779* (immunoglobulin superfamily member 1)
311 was linked with six carcass composition traits, and *ENSGALG00000050638*
312 (immunoglobulin-like receptor CHIR2D-878 precursor) was associated with shank

313 length (SL12) (Supplementary Table S8).

314 As expected, we also found that the promoter region of *IGF2BP1* was associated
315 with growth traits, including claw weight (CW1), the ratio of claw weight to body
316 weight (CR), double pinion weight (DPW) and semi-evisceration weight (SEW), based
317 on PAV-based GWAS of both 1-2 kb and 2-3 kb upstream regions (Figure 3e-f;
318 Supplementary Table S8). The most significant association was identified between
319 *IGF2BP1* and CW1 ($p = 1.92E-07$) based on 1-2 kb PAV-based GWAS (Figure 3h-i).

320 **Dissection of the structure and function of *IGF2BP1* promoter region**

321 To dissect the structure of *IGF2BP1* promoter region, we comprehensively analyzed
322 the results of PCR and WGS read mapping. Three alleles of *IGF2BP1* promoter region
323 were identified, which were defined as W (wild type), L1 allele (3.2 kb deletion at
324 GRCg6a chr27:6082202-6085435), and L2 allele (1.5 kb deletion at chr27:6083984-
325 6085538) (Figure 4a; Supplementary Figure S9a). We conducted allele specific PCR to
326 genotype these three alleles in wild, native and commercial chickens (Figure 4a-b). As
327 expected, the W allele was dominant in native breeds and wild relatives. In contrast, all
328 the two commercial broiler breeds and commercial crossed chickens mainly had
329 absence variant (L1 or L2 alleles). Absence variant (L1 or L2) was also dominant in
330 commercial layer breeds, except for White Leghorn breeds. This result is consistent
331 with the distribution of 1-2 kb and 2-3 kb upstream region PAV frequency for the
332 *IGF2BP1* promoter region, which showed that commercial breeds were almost uniform
333 for the mutant absence variant. We also compared the promoter region of *IGF2BP1*
334 between high weight selection (HWS) lines and low weight selections (LWS) line using
335 their pool sequencing data (Lillie, et al. 2018). We found that L1 was fixed in high body
336 weight lines, whereas W was fixed in the low body weight lines, including the relaxed
337 selection lines (Supplementary Figure S9b). It implies that W and L1 alleles have been
338 selected to be fixed at an earlier time, before the divergence of relaxed selection lines.

339 Via the single genotype marker association analysis, the associations between the
340 L1 allele and the body size, body weight or carcass composition still hold true when
341 enlarging the sample size of Gushi×Anak F2 population to 734 (Figure 5;

342 Supplementary Figure S10). The associated traits include claw weight (CW1, CR),
343 shank length (SL12), breast bone length (BBL8 and BBL12), wing weight (DPW),
344 evisceration weight (EW and SEW), head weight (HW1), carcass weight (CW), leg
345 weight (LW and LMW), pelvis breadth (PB12), shank girth (SG12 and SG8), body
346 slanting length (BSL8 and BSL12), gizzard weight (GW), body weight (BWHR, BW6
347 and BW10), and growth rate (GR0_4). In those traits, the L1L1 genotype was always
348 linked to better performance of production (such as body size, carcass weight and body
349 weight) than the WW genotype (Figure 5; Supplementary Figure S10). The significant
350 association between the *IGF2BP1* genotype and body size confirms the PAV-based
351 GWAS results in promoter regions (Figure 3). Of these, associations are most
352 significant in traits CW1 ($p = 2.32E-14$) and CR ($p = 3.70E-12$), which account for
353 4.01% and 3.85% of the phenotypic variations, respectively. Interestingly, we observed
354 a larger effect of L1 in females relative to males, which explained 11.5% in females and
355 7.3% in males of the phenotypic variations for CW1 trait (Figure 5a). Besides, a female
356 phenotype variation of 8.2% for DPW and 6.2% for SL12 was explained by L1. These
357 associations are also consistent with the chicken and duck SNP-based GWAS results,
358 which indicated that SNPs located near *IGF2BP1* were associated with body weight,
359 head weight, gizzard weight, chest width, leg weight, and wing weight (Ma, et al. 2018;
360 Zhou, et al. 2018; Wang, Bu, et al. 2020; Wang, Cao, et al. 2020; Zhang, et al. 2020).
361 Unexpectedly, L2 allele was not found in the F2 population.

362 To further verify the molecular effects of the deletions, luciferase expression levels
363 were investigated to represent the transcriptional activity through transfecting three
364 kinds of recombinant plasmids (pGL3-L1, pGL3-L2, and pGL3-W) into chicken DF-1
365 cells (Figure 6a). Before performing the luciferase activity experiment, we screened the
366 genome region which inserted the pGL3 construct, and confirmed that we did not find
367 any difference except the L1 and L2 deletion. Therefore, the activity difference among
368 the three constructs was derived from the deletions. The transcriptional activities of
369 these two deletions (L1 and L2) were significantly higher than that of wild type (W).
370 Further, the activity of the L1 genotype is also higher than that of L2 (Figure 6a).

371 Subsequently, we compared the mRNA expression level between the L1L1 (Ross 308)
372 and WW genotypes (Gushi chicken). The expression of *IGF2BP1* mRNA in WW
373 genotype is significantly lower than that in the L1L1 genotype in almost all investigated
374 tissues at 6 weeks of age (Figure 6b). To reduce the difference in the genetic background
375 among individuals with different genotypes and investigate the effect of deletions more
376 accurately, we performed cross-breeding between chickens with the same heterozygous
377 genotypes (L1WxL1W and L2WxL2W) to generate L1L1, L2L2 and WW genotype
378 chicken with half-sib or full-sib relationship, and then compared the expressions of
379 *IGF2BP1*. In spleen and duodenum tissues at 3 weeks of age, we observed higher
380 expressions in L1L1 and L2L2 than WW genotype, while L1L1 also showed a higher
381 value than L2L2 (Figure 6c). This is completely consistent with the result of
382 transcriptional activity (Figure 6a). We also observed three conserved elements located
383 in L1 deletion based on 77 vertebrates basewise PhyloP conservation score
384 (<https://hgdownload.soe.ucsc.edu/goldenPath/galGal6/phastCons77way/>), of which
385 one conserved element located in L2 deletion (Figure 4a). These suggest the functional
386 importance of three conserved elements which possibly regulating *IGF2BP1*
387 expression.

388 **Investigation of the genomic regions flanking the deletion**

389 Since *IGF2BP1* was also reported as the body size candidate gene using SNP-based
390 studies (Sheng, et al. 2013; Ma, et al. 2018; Wang, Bu, et al. 2020; Wang, Cao, et al.
391 2020; Zhang, et al. 2020), we explored the SNPs within the region from 10 kb upstream
392 of L1 deletion and 10 kb downstream in order to test if any of the signal driving SNPs
393 in previous studies could be the causal. SNP calling was done using the same 664
394 individual sequencing data for building the pan-genome; however, 210 individuals were
395 excluded for the low quality of SNP calling or their heterozygous genotype. The
396 remaining individuals include 325 WW, 117 L1L1, and 12 L2L2 samples. We searched
397 for SNPs that associated with the deletions in three different ways, L1L1 vs. WW, L2L2
398 vs. WW, and L1L1 + L2L2 vs. WW. Altogether, five associated SNPs were detected.
399 Among them, the highest PhyloP conservation score is 0.97, and that SNP (chr27:

400 6087849) is not within a conservation element. The other four SNPs have negative
401 conservation scores. This implies that none of these five associated SNPs are highly
402 conserved, which supports that the deletion is likely to be the only functional mutation
403 within this region.

404 **Discussion**

405 **Construction of the first chicken pan-genome and dissection of genetic changes in** 406 **the chicken population**

407 Here, we constructed the first pan-genome of chicken, capturing ~66.5 Mb novel
408 sequences that are absent from the reference genome (GRCg6a). Similar novel
409 additional pan-genome sequences were captured in pig (Tian, et al. 2020) (~72.5 Mb),
410 human (Sherman, et al. 2019) (~296.5 Mb), and plants (Yao, et al. 2015; Golicz, Bayer,
411 et al. 2016; Montenegro, et al. 2017) (15.8 Mb ~ 350 Mb). Absent sequences from the
412 reference genome were predicted to encode additional 4,063 high-confidence genes.
413 We also found about one-third of the gene PAV is variable among the 268 individuals
414 used for PAV calling. This highlights the heterogeneity of genetic makeup among
415 chicken breeds and shows a potential utility for further breeding (Gao, et al. 2019).

416 We observed that red jungle fowls and native chickens contained most of the genetic
417 diversity of chickens, while limited genetic diversity was found in commercial chickens
418 (Figure 2). This result is consistent with known reductions in genetic diversity of
419 modern livestock compared to their wild ancestors (Malomane, et al. 2019; Frantz, et
420 al. 2020). Similarly, peach (Guo, et al. 2020), chickpea (Varshney, et al. 2019) and
421 tomato (Gao, et al. 2019) pan-genome studies found that their wild relatives and
422 landraces are more genetically diverse compared with modern cultivars. We also found
423 that intra-species gene content variation can be affected by selection, genetic drift, or
424 hybridization (Supplementary Figure S3-S5). We proposed that the reduction of genetic
425 diversity in commercial chickens might occur due to intensive artificial selection during
426 breeding, but other factors can not be ruled out, such as genetic drift.

427 **PAVs are associated with physiological traits and the presence frequency of** 428 **immune related loci was reduced during modern chicken breeding**

429 We found that the promoter region PAV of genes showed associations with
430 physiological related traits, such as *LDLR* and olfactory receptors (Supplementary
431 Table S8). Lipid accumulation can enhance *LDLR* expression leading to an increase of
432 serum creatinine (Sun, et al. 2013; Zhang, et al. 2016). *LDLR* knockout mouse and rat
433 showed substantial increases in plasma creatinine (Bisgaard, et al. 2016; Sithu, et al.
434 2017). Variation in the promoter region of *LDLR* may reduce its expression and further
435 upregulate serum creatinine level. Olfactory receptors were first discovered in the
436 olfactory epithelium, functioning in odorant recognition involving various
437 physiological behaviors, such as food choice and intake. However, recent studies
438 indicate that these genes are also expressed in the intestinal tract (Priori, et al. 2015;
439 Kim, et al. 2017; Kotlo, et al. 2020), and olfactory receptors play a role in intestinal
440 inflammatory reaction (Kotlo, et al. 2020), secretion (Kim, et al. 2017) and microbiota
441 metabolites (Priori, et al. 2015). We also found olfactory receptor 14C36-like gene
442 associated with ileum length (IL), jejunum length (JL) and cecum length (CL)
443 (Supplementary Table S8). It is thus possible that olfactory receptors are involved in
444 feed digestion and conversion via regulation of intestine development and thus were
445 under selection during modern breeding (Supplementary Figure S7a).

446 We observed the presence frequency of immune related gene or promoter region
447 (including MHC and immunoglobulin) decreased in commercial chicken compared
448 with the native breed. Of these, some immune gene PAVs showed significant
449 association with production traits (Supplementary Table S8). This is consistent with a
450 previous report that a high immune response is negatively correlated with chicken egg
451 production and body weight (Warner, et al. 1987). MHC genes are involved in immune
452 recognition and susceptibility to infectious disease (Sommer 2005). There is a possible
453 genetic linkage between MHC genes and growth or reproduction genes (Warner, et al.
454 1987). Another possible explanation is that increased productivity may also increase
455 the metabolic burden of immune gene maintenance in modern breeds. A trade-off might
456 occur between the conservation of production-related genes and the loss of immune-
457 related genes due to human selection for desirable production traits (van der Most, et

458 al. 2011).

459 ***IGF2BP1* deletion is the causal variant for a major QTL associated with body size**

460 Many QTGs or QTLs associated with chicken growth traits have been identified, of
461 which loci located at chromosomes 27, 4 and 1 have the largest impact on growth in
462 chicken (Sheng, et al. 2013; Ma, et al. 2018; Wang, Bu, et al. 2020; Wang, Cao, et al.
463 2020; Zhang, et al. 2020). To our knowledge, the study in 2003 was the first time to
464 report that a large QTL region located between 4.0 Mb and 6.1 Mb in chromosome 27
465 was associated with chicken body size (Kerje, et al. 2003). After that, many studies
466 identified the chicken growth trait QTL in chromosome 27, including the gene
467 *IGF2BP1* by SNP-based GWAS (Sheng, et al. 2013; Ma, et al. 2018; Wang, Bu, et al.
468 2020; Wang, Cao, et al. 2020). Our previous GWAS also revealed a signal peak
469 correlated to body size trait, which was located at the genomic upstream of *IGF2BP1*
470 (Zhang, et al. 2020). GWAS in duck also revealed SNPs located at the genomic
471 upstream region of *IGF2BP1* that showed significant association with body size traits,
472 while a higher expression level of *IGF2BP1* is correlated to better performance (Zhou,
473 et al. 2018). Altogether, *IGF2BP1* is a potential major gene associated with body size
474 traits, but the causal variant regulating these traits has not been reported previously.

475 In this study, using a genotype-phenotype association, we found two mutant alleles
476 in the *IGF2BP1* promoter region that contributed to larger body size. We also observed
477 a stronger association in females than males (Figure 5; Supplementary Figure S10). We
478 compared the phenotypes among L1W, L1L1, and WW chickens to estimate the
479 inheritance mode of the deletions. Taking the CW1 trait as an example, we found no
480 significant difference in CW1 between L1W and L1L1 ($p = 0.68$) in males, while both
481 are significantly heavier than WW (L1W vs WW, $p = 1.26E-3$, L1L1 vs WW, $p = 2.60E-$
482 5). We inferred that there is a possible dominant effect of L1 against W in males. In
483 females, however, we found no significant difference between WW and L1W ($p = 0.42$),
484 while L1L1 are significantly heavier than L1W ($p = 5.35E-7$) and WW ($p = 4.0E-6$).
485 There is a possible recessive effect of L1 against W in females. One possible reason is
486 that this autosomal deletion locus shows sex-influenced inheritance, with a dominant

487 effect in males and a recessive effect in females. There may be a putative binding site
488 of androgen-mediated transcription factor located on this deletion region. We also found
489 three conserved elements based on 77 vertebrates basewise PhyloP conservation score,
490 suggesting a putative regulatory function (Figure 4a). These deletions in the promoter
491 region may increase *IGF2BP1* expression by upregulating its transcriptional activity
492 (Figure 6). Further studies are required to elucidate the upstream regulatory pathway.

493 Together with our GWAS analysis, the mutant genotype is associated with higher
494 expression of *IGF2BP1* and improved productivity traits (Figure 5; Figure 6). Our
495 findings are consistent with findings that higher expression of *IGF2BP1* is linked to the
496 larger body size in duck (Zhou, et al. 2018). Although the *IGF2BP1* mutation only
497 explains a moderate 2-4% of phenotypic variation, this is in fact a substantial effect for
498 a complex quantitative trait like body size. For instance, in humans two key variants for
499 lean body mass explained 0.23% and 0.16% of the variance (Zillikens, et al. 2017) and
500 ~50 variants for height only explain ~5% of the variance (Yang, et al. 2010). After
501 examining the flanking regions of the deletion, the only five SNPs correlated to the
502 deletion showed extremely low conservation scores implying that the deletion is the
503 unique functional variant in this region. Based on this combined evidence, we propose
504 that the deletion in the *IGF2BP1* promoter region is the causal variant for the QTL
505 located at chromosome 27 that was previously reported to be related to body size in
506 chicken.

507 **Conclusion**

508 Collectively, this first chicken pan-genome provides a foundation for future chicken
509 population genetics and evolutionary genomics studies. PAV analysis offers an
510 opportunity to uncover genomic architecture and identify the change of gene content
511 during domestication and improvement, helping the designing of future chicken breeds
512 with desired traits. We dissect the causal variant of one of the major QTL contributing
513 to body size in chicken using PAV-based GWAS. The deletions that we found can be
514 applied as markers for breeding programs using marker-assisted selection. As pan-
515 genomic studies become more common, PAV-based GWAS will provide a powerful

516 complement to SNP-based GWAS for identifying functional variants of economically
517 or evolutionary important traits.

518 **Materials and Methods**

519 **Genomic sequencing of chicken**

520 A total of 868 individuals were used in this study, of which 664 were used to construct
521 the chicken pan-genome (Supplementary Table S1). We downloaded 509 accessions,
522 published in recent genome resequencing studies (Fan, et al. 2013; Wang, et al. 2015;
523 Ulfah, et al. 2016; Li, Che, et al. 2017; Lawal, et al. 2018; Qanbari, et al. 2019; Huang,
524 et al. 2020; Wang, Thakur, et al. 2020b), from the National Center for Biotechnology
525 Information (NCBI) Sequence Read Archive database (Supplementary Table S1).
526 Sequencing data of 150 Henan indigenous chickens and 204 Gushi×Anak F2
527 individuals were generated in this study, and further data for an additional 5 Xichuan
528 black-bone chickens were generated in our previous study (Li, Sun, et al. 2020).
529 Genomic DNA was extracted from chicken blood using Qiagen DNeasy Kit. Paired-
530 end libraries with ~500 bp insert size were constructed and then subjected to sequencing
531 using the BGISEQ-500 platform to generate paired-end 150 bp reads (BGI Genomics
532 Co., Ltd. and Beijing Fuyu Biotechnology Co., Ltd, China). We also downloaded 10
533 pool sequencing data, including 5 HWS and 5 LWS pool data from the NCBI database
534 using project number PRJNA516366 (Lillie, et al. 2018).

535 **Pan-genome construction and annotation**

536 Raw reads were processed to remove low quality reads and generate adaptor free clean
537 reads using Trimmomatic (v0.36) (Bolger, et al. 2014). The pan-genome was
538 constructed by a reference based iterative mapping and assembly approach using the
539 GRCg6a assembly as a starting reference genome (Golicz, Batley, et al. 2016; Golicz,
540 Bayer, et al. 2016). The reference-based iterative mapping and assembly approach
541 (Golicz, Batley, et al. 2016; Golicz, Bayer, et al. 2016) was first applied in a pan-
542 genome study of the crop *B.oleracea*. This approach allows using sequencing data from
543 a large range of individuals from different populations to construct a pan-genome.
544 Briefly, clean reads were mapped to the reference genome (Ensemble
545 *Gallus_gallus.GRCg6a.dna.toplevel.fa*) using bowtie2 (v2.3.5.1) (Langmead and
546 Salzberg 2012). Unmapped reads were extracted using SAMtools and then assembled

547 using MaSuRCA v3.3.1 (Zimin, et al. 2013). After pan-genome construction, newly
548 assembled contigs of non-reference sequences with length larger than 500 bp were kept.
549 Contaminant sequences were filtered by the following two steps. Firstly, contigs were
550 aligned using blastn v2.9.0 (Camacho, et al. 2009) against the NT database (v5, 07-03-
551 2019) of contaminant taxid groups, which includes archaea, viruses, bacteria, fungi and
552 Viridiplantae. Secondly, the remaining contigs were classified and filtered using
553 Kraken2 (v 2.0.9-beta) using the kraken2-microbial database, which consists of archaea,
554 bacteria, fungi, protozoa, viral and human sequences.
555 (https://lomanlab.github.io/mockcommunity/mc_databases.html) (Wood, et al. 2019).
556 The unclassified contigs were defined as contamination-free. The final contamination-
557 free non-reference sequences and the reference Gallus/GRCg6a genome were merged
558 to generate the chicken pan-genome.

559 A custom repeat library was constructed by scanning the final non-reference
560 sequence using RepeatModeler (v1.0.11) (Flynn, et al. 2020). A custom repeat library
561 and the RepBase database (downloaded in June 2019) of vertebrates were used to detect
562 the repeat sequences with RepeatMasker (v4.0.8) (Tarailo-Graovac and Chen 2009).
563 The MAKER2 annotation pipeline was used to obtain a set of high-confidence
564 annotation based on RNA-seq evidence, homologous protein evidence and *ab initio*
565 gene prediction evidence (Holt and Yandell 2011). RNA-seq evidence was generated
566 using Hisat2-Stringtie pipeline (Pertea, et al. 2016) with published data from available
567 tissues (Supplementary Table S2). Protein sequences of chicken, human and other
568 mammals and vertebrates were collected from the Uniprot database
569 (<https://www.uniprot.org/>). *Ab initio* gene prediction was implemented using SNAP
570 (Korf 2004) and Augustus (Stanke, et al. 2006) with the ‘chicken’ model selected.
571 Finally, redundant assembled protein sequences were filtered with CD-HIT (Fu, et al.
572 2012) (-c 0.9 -n 5 -M 16000 -T 18) with the threshold of 90% similarity.

573 **PAV calling**

574 Gene PAV was determined based on the cumulative coverage of exons of each gene.
575 The longest transcripts were retrieved as the gene body to avoid redundant gene counts.

576 If at least two reads covered more than 5% cumulative coverage of all exons, this gene
577 was defined as present in an individual. Otherwise, it was defined as absent (Golicz,
578 Bayer, et al. 2016). Clean reads were aligned to the pan-genome using BWA-MEM
579 (v0.7.17) (Li and Durbin 2009) with default parameters and the sequences depth of each
580 sample was captured using Mosdepth package (v0.2.5)(Pedersen and Quinlan 2018).
581 High-depth sequencing data (>30x) is preferable to increase the robustness of PAV
582 analysis; however, it is not economical to sequence large samples numbers at this depth.
583 Low-depth data (<15x) is a viable and more economical means to carry out PAV
584 analysis in large sets of diverse samples (Gao, et al. 2019; Sherman, et al. 2019;
585 Jayakodi, et al. 2020). To estimate the impact of the sequencing depth on gene PAV
586 calling, we extracted reads from reference genome individual with varying depths of
587 sequences to determine the minimum sequence depth required to call a confident gene
588 PAV. An average sequence depth of 10x was considered as the threshold for including
589 a sample since this threshold is estimated to allow a 99.94% recovery rate of gene PAV
590 (Gao, et al. 2019) (Supplementary Figure S2a). We also performed additional
591 simulation analysis using sequencing data of random seven breeds and found
592 98.4%~99.5% of pan-genome genes can be called when the average sequence depth
593 reaching 10x (Supplementary Figure S2b). Thus, to get a high confident PAV matrix,
594 individuals with an average depth above 10x were kept to perform gene PAV calling.
595 Additionally, the sequencing data of red jungle fowls in Thailand were reported to be
596 contaminated by domestic chicken sequences and were removed for the PAV calling
597 (Qanbari, et al. 2019; Wang, Thakur, et al. 2020b).

598 PAV calling for promoter region was performed using the same method of gene
599 PAV calling that is described above but based on the gene promoter regions. We divided
600 the promoter region into three 1 kb windows based on the distance to the transcription
601 start site. The three blocks were in the 0-1 kb, 1-2 kb and 2-3 kb regions upstream to
602 the transcription start site of genes in the reference genome. A PAV was considered as
603 present if more than 50% cumulative coverage with at least two reads was identified,
604 otherwise, it was considered absent (Golicz, Bayer, et al. 2016).

605 **PAV analysis**

606 The gene PAV matrix was subjected to population genetic analysis. Principal
607 component analysis and neighbour-joining phylogenetic analysis were conducted using
608 TASSEL5 (Bradbury, et al. 2007). To identify the PAV with frequency significantly
609 changed during improvement, the PAV frequency of each gene was compared between
610 the native breeds and commercial breeds. Fisher's exact test was employed to identify
611 significant PAV with false discovery rate (FDR) 0.001 (Gao, et al. 2019). Significantly
612 increased genes were defined as genes having a significantly higher frequency in the
613 commercial breeds than the native breeds. Inversely, we consider genes with a
614 significantly lower frequency as significantly decreased genes. To identify the promoter
615 region with significantly changed during chicken improvement, PAV patterns were also
616 analyzed using the same method as gene PAV frequency calculation.

617 **PAV-based GWAS**

618 PAV-based genome-wide association study (GWAS) was also implemented to identify
619 the candidate genes associated with 151 traits in a GushixAnak F2 mapping population
620 with 204 individuals. To reduce bias, gene PAVs were removed if they were located on
621 sex chromosomes or showed a minor allele frequency less than 0.05. A general linear
622 model (GLM) was employed for association analysis using TASSEL5 (Bradbury, et al.
623 2007), with sex and the first five PCA eigenvectors defined as fixed effects. A
624 Bonferroni test was used to define the genome-wide significant (0.05/number of loci)
625 or suggestive (0.1/number of loci) cut-off threshold.

626 **GO annotation**

627 Functional annotation of the pangenome was performed using command line Blast2GO
628 (Conesa et al., 2005) v2.5. The pan-genome genes were aligned to the proteins in the
629 Uniref90 database (downloaded on Sep 2019) using BLASTP (Camacho et al., 2009),
630 and only alignments with E-values $< 1 \times 10^{-5}$ were used. Then, the BLAST results were
631 reformatted to satisfy Blast2GO naming requirements. Gene ontology annotation of the
632 variable genes was conducted by the R package topGO (Alexa, et al. 2006) using
633 Fisher's exact test with the approach 'elim' used to correct for multiple comparisons.

634 **Genotyping of *IGF2BP1* PAV and association analysis**

635 Three primers, including one forward and two reverse primers, were designed based on
 636 the sequence of the *IGF2BP1* promoter region (Figure 4a). One pair of primer, *Asp-F*
 637 and *Asp-R*, was used for genotyping L1 and W allele, while another pair, *2k-F* and *Asp-*
 638 *R*, was used to genotype L2 and W (Figure 4a; Supplementary Table S8). PCR reaction
 639 was conducted as described below: 5 pmol of each primer, 100 ng of genomic DNA, 2
 640 ul 10 x PCR buffer (Takara), 100 uM dNTP mixture and 1 ul Taq polymerase.
 641 Association analysis of the validation population was conducted between genotypes
 642 (L1L1, L1W and WW) in *IGF2BP1* PAV and 151 traits (Supplementary Note) in F2
 643 population with 734 individuals using GLM as described as above. The value of marker
 644 R squared was used to explain the phenotype variation of *IGF2BP* loci, as computed
 645 from the marker sum of squares (SS) after fitting all other model terms divided by the
 646 total SS (Bradbury, et al. 2007).

647 **Functional assay of *IGF2BP1* promoter region and *IGF2BP1* expression**

648 Three kinds of *IGF2BP1* promoter region (L1, L2 and W) were cloned into pGL3-Basic
 649 luciferase vector (Promega) using *Clone-F* and *Clone-R* primer (Supplementary Table
 650 S9). All recombinant plasmids, together with the pRL-TK plasmid (Promega), were
 651 transfected into DF-1 (chicken fibroblast cell) cell line. After 48 hours, the
 652 transcriptional activity was investigated by the Dual-Luciferase Reporter Assay System
 653 (Promega). Quantitative PCR was conducted to investigate the mRNA level of
 654 *IGF2BP1* using primer *IGF2BP1-qF* and *IGF2BP1-qR* (Supplementary Table S9). The
 655 relative expression level of *IGF2BP1* was normalized by GAPDH using the $2^{-\Delta\Delta ct}$
 656 method.

657 **Investigating the flanking SNPs of deletions**

658 The deletion and its flanking regions (chr27:6072202-6095435) were analyzed by
 659 GATK (v3.8) pipeline (McKenna, et al. 2010) using the same 664 individuals for
 660 building the pan-genome. Genotypes of the *IGF2BP1* deletion of each sample were
 661 determined by the GATK results and manually checking the alignments by IGV
 662 (version 2.4.3). Then samples with the same genotypes were grouped together, and the

663 SNP associated with the *IGF2BP1* genotypes was defined as 1) significant in the Chi-
664 squared test, 2) the mutant allele of the SNP has an allelic frequency higher than 0.8 in
665 the deletion group, and 3) the allelic frequency difference between the two compared
666 groups greater than 0.5. Three different deletion groups were used in three different
667 comparisons. They are L1L1 group, L2L2 group, and L1L1 + L2L2 group, all compared
668 with WW group, respectively.

669 **Availability of Data and Materials**

670 All the sequence data generated in this study have been deposited in the National
671 Genomics Data center (<https://bigd.big.ac.cn>) with the accession codes PRJCA004227
672 and PRJCA004441. Downloaded sequence data used in this study were presented in
673 Supplementary Table S1. The chicken pan-genome and relevant data are available at
674 the DRYAD database (<https://doi.org/10.5061/dryad.7pvmcvds1>)

675 **Author contribution**

676 K.W., H.H. and W. L. designed analysis, performed analysis and wrote manuscript;
677 W.L., C.Z., Y.L., J.W., L.Y., and X.F. performed the wet-lab experiment; X.K., Y.T.,
678 G.S., D.L., Y.Z., R.H., R.J., F.Y., Y.W., Z.L., G.L, X.L. contributed to sample collection
679 and construction of F2 resource population. J. L. and A.S. assisted with data analysis
680 and manuscript revision. W. L., D.E. and X. K. conceived research designed analysis
681 and revised manuscript.

682 **Acknowledgements and funding information**

683 We thank Leif Andersson for comments on an earlier version of this manuscript. We
684 thank Longxian Zhang and Jiangying Huang for help on computing resource. This work
685 was supported by the Program for Innovation Research Team of the Ministry of
686 Education (IRT16R23), National Natural Science Foundation of China (31902144) and
687 the Scientific Studio of Zhongyuan Scholars (30601985). H.H. thanks the China
688 Scholarship Council for supporting his studies at the University of Western Australia.

689 **Ethics declarations**

690 Ethics approval for this study was obtained from Henan Agricultural University.

691 **Competing interests**

692 The authors declare that they have no competing interests.

693 **References**

- 694 Abe H, Gemmell NJ. 2014. Abundance, arrangement, and function of sequence motifs in the chicken
695 promoters. *BMC Genomics* 15:900.
- 696 Alexa A, Rahnenfuhrer J, Lengauer T. 2006. Improved scoring of functional groups from gene expression
697 data by decorrelating GO graph structure. *Bioinformatics* 22:1600-1607.
- 698 Barton NH. 2001. The role of hybridization in evolution. *Mol Ecol* 10:551-568.
- 699 Bayer PE, Golicz AA, Scheben A, Batley J, Edwards D. 2020. Plant pan-genomes are the new reference.
700 *Nature Plants* 6:914-920.
- 701 Bell JL, Wachter K, Muhleck B, Pazaitis N, Kohn M, Lederer M, Huttelmaier S. 2013. Insulin-like
702 growth factor 2 mRNA-binding proteins (IGF2BPs): post-transcriptional drivers of cancer progression?
703 *Cell Mol Life Sci* 70:2657-2675.
- 704 Bertolotti AC, Layer RM, Gundappa MK, Gallagher MD, Pehlivanoglu E, Nome T, Robledo D, Kent
705 MP, Rosaeg LL, Holen MM, et al. 2020. The structural variation landscape in 492 Atlantic salmon
706 genomes. *Nature Communications* 11.
- 707 Bickhart DM, Liu GE. 2014. The challenges and importance of structural variation detection in livestock.
708 *Frontiers in Genetics* 5.
- 709 Bisgaard LS, Bosteen MH, Fink LN, Sorensen CM, Rosendahl A, Mogensen CK, Rasmussen SE, Rolin
710 B, Nielsen LB, Pedersen TX. 2016. Liraglutide Reduces Both Atherosclerosis and Kidney Inflammation
711 in Moderately Uremic LDLr^{-/-} Mice. *PLoS One* 11:e0168396.
- 712 Bolger AM, Lohse M, Usadel B. 2014. Trimmomatic: a flexible trimmer for Illumina sequence data.
713 *Bioinformatics* 30:2114-2120.
- 714 Bradbury PJ, Zhang Z, Kroon DE, Casstevens TM, Ramdoss Y, Buckler ES. 2007. TASSEL: software
715 for association mapping of complex traits in diverse samples. *Bioinformatics* 23:2633-2635.
- 716 Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL. 2009. BLAST+:
717 architecture and applications. *BMC Bioinformatics* 10:421.
- 718 DANFORTH CH. 1958. GALLUS SONNERATI AND THE DOMESTIC FOWL. *Journal of Heredity*
719 49:167-170.
- 720 DESTA TT. 2019. Phenotypic characteristic of junglefowl and chicken. *World's poultry science journal*
721 2019 v.75 no.1:pp. 69-82.
- 722 Dorshorst B, Molin AM, Rubin CJ, Johansson AM, Stromstedt L, Pham MH, Chen CF, Hallbook F,
723 Ashwell C, Andersson L. 2011. A complex genomic rearrangement involving the endothelin 3 locus
724 causes dermal hyperpigmentation in the chicken. *PLoS Genet* 7:e1002412.
- 725 Duan Z, Qiao Y, Lu J, Lu H, Zhang W, Yan F, Sun C, Hu Z, Zhang Z, Li G, et al. 2019. HUPAN: a pan-
726 genome analysis pipeline for human genomes. *Genome Biol* 20:149.
- 727 Edwards AW. 2008. G. H. Hardy (1908) and Hardy-Weinberg equilibrium. *Genetics* 179:1143-1150.
- 728 Eriksson J, Larson G, Gunnarsson U, Bed'hom B, Tixier-Boichard M, Stromstedt L, Wright D, Jungerius
729 A, Vereijken A, Randi E, et al. 2008. Identification of the yellow skin gene reveals a hybrid origin of the
730 domestic chicken. *PLoS Genet* 4:e1000010.
- 731 Fallahshahroudi A, Sorato E, Altimiras J, Jensen P. 2019. The Domestic BCO2 Allele Buffers Low-
732 Carotenoid Diets in Chickens: Possible Fitness Increase Through Species Hybridization. *Genetics*
733 212:1445-1452.
- 734 Fan WL, Ng CS, Chen CF, Lu MY, Chen YH, Liu CJ, Wu SM, Chen CK, Chen JJ, Mao CT, et al. 2013.

735 Genome-wide patterns of genetic variation in two domestic chickens. *Genome Biol Evol* 5:1376-1392.
736 Farre D, Bellora N, Mularoni L, Messeguer X, Alba MM. 2007. Housekeeping genes tend to show
737 reduced upstream sequence conservation. *Genome Biol* 8:R140.
738 Flynn JM, Hubley R, Goubert C, Rosen J, Clark AG, Feschotte C, Smit AF. 2020. RepeatModeler2 for
739 automated genomic discovery of transposable element families. *Proc Natl Acad Sci U S A* 117:9451-
740 9457.
741 Frantz LAF, Bradley DG, Larson G, Orlando L. 2020. Animal domestication in the era of ancient
742 genomics. *Nat Rev Genet* 21:449-460.
743 Fu L, Niu B, Zhu Z, Wu S, Li W. 2012. CD-HIT: accelerated for clustering the next-generation
744 sequencing data. *Bioinformatics* 28:3150-3152.
745 Furman C, Sieminski AL, Kwiatkowski AV, Rubinson DA, Vasile E, Bronson RT, Fassler R, Gertler FB.
746 2007. Ena/VASP is required for endothelial barrier function in vivo. *Journal of Cell Biology* 179:761-
747 775.
748 Gao L, Gonda I, Sun H, Ma Q, Bao K, Tieman DM, Burzynski-Chang EA, Fish TL, Stromberg KA,
749 Sacks GL, et al. 2019. The tomato pan-genome uncovers new genes and a rare allele regulating fruit
750 flavor. *Nat Genet* 51:1044-1051.
751 Gerdol M, Moreira R, Cruz F, Gomez-Garrido J, Vlasova A, Rosani U, Venier P, Naranjo-Ortiz MA,
752 Murgarella M, Greco S, et al. 2020. Massive gene presence-absence variation shapes an open pan-
753 genome in the Mediterranean mussel. *Genome Biology* 21.
754 Golicz AA, Batley J, Edwards D. 2016. Towards plant pangenomics. *Plant Biotechnol J* 14:1099-1105.
755 Golicz AA, Bayer PE, Barker GC, Edger PP, Kim H, Martinez PA, Chan CK, Severn-Ellis A, McCombie
756 WR, Parkin IA, et al. 2016. The pangenome of an agronomically important crop plant *Brassica oleracea*.
757 *Nat Commun* 7:13390.
758 Guo J, Cao K, Deng C, Li Y, Zhu G, Fang W, Chen C, Wang X, Wu J, Guan L, et al. 2020. An integrated
759 peach genome structural variation map uncovers genes associated with fruit traits. *Genome Biol* 21:258.
760 Guo Y, Gu X, Sheng Z, Wang Y, Luo C, Liu R, Qu H, Shu D, Wen J, Crooijmans RP, et al. 2016. A
761 Complex Structural Variation on Chromosome 27 Leads to the Ectopic Expression of HOXB8 and the
762 Muffs and Beard Phenotype in Chickens. *PLoS Genet* 12:e1006071.
763 Hansen TV, Hammer NA, Nielsen J, Madsen M, Dalbaeck C, Wewer UM, Christiansen J, Nielsen FC.
764 2004. Dwarfism and impaired gut development in insulin-like growth factor II mRNA-binding protein
765 1-deficient mice. *Mol Cell Biol* 24:4448-4464.
766 Henrich M, Buckler KJ. 2013. Cytosolic calcium regulation in rat afferent vagal neurons during anoxia.
767 *Cell Calcium* 54:416-427.
768 Hillier LW, Miller W, Birney E, Warren W, Hardison RC, Ponting CP, Bork P, Burt DW, Groenen MAM,
769 Delany ME, et al. 2004. Sequence and comparative analysis of the chicken genome provide unique
770 perspectives on vertebrate evolution. *Nature* 432:695-716.
771 Holt C, Yandell M. 2011. MAKER2: an annotation pipeline and genome-database management tool for
772 second-generation genome projects. *BMC Bioinformatics* 12:491.
773 Huang H, Weng H, Sun W, Qin X, Shi H, Wu H, Zhao BS, Mesquita A, Liu C, Yuan CL, et al. 2018.
774 Recognition of RNA N(6)-methyladenosine by IGF2BP proteins enhances mRNA stability and
775 translation. *Nat Cell Biol* 20:285-295.
776 Huang X, Otecko NO, Peng M, Weng Z, Li W, Chen J, Zhong M, Zhong F, Jin S, Geng Z, et al. 2020.
777 Genome-wide genetic structure and selection signatures for color in 10 traditional Chinese yellow-
778 feathered chicken breeds. *BMC Genomics* 21:316.

779 Imsland F, Feng C, Boije H, Bed'hom B, Fillon V, Dorshorst B, Rubin CJ, Liu R, Gao Y, Gu X, et al.
780 2012. The Rose-comb mutation in chickens constitutes a structural rearrangement causing both altered
781 comb morphology and defective sperm motility. *PLoS Genet* 8:e1002775.
782 Jayakodi M, Padmarasu S, Haberer G, Bonthala VS, Gundlach H, Monat C, Lux T, Kamal N, Lang D,
783 Himmelbach A, et al. 2020. The barley pan-genome reveals the hidden legacy of mutation breeding.
784 *Nature*.
785 Jing Z, Wang X, Cheng Y, Wei C, Hou D, Li T, Li W, Han R, Li H, Sun G, et al. 2020. Detection of CNV
786 in the SH3RF2 gene and its effects on growth and carcass traits in chickens. *BMC Genet* 21:22.
787 Keightley PD, Lercher MJ, Eyre-Walker A. 2005. Evidence for widespread degradation of gene control
788 regions in hominid genomes. *PLoS Biol* 3:e42.
789 Kerje S, Carlborg O, Jacobsson L, Schutz K, Hartmann C, Jensen P, Andersson L. 2003. The twofold
790 difference in adult size between the red junglefowl and White Leghorn chickens is largely explained by
791 a limited number of QTLs. *Anim Genet* 34:264-274.
792 Kerstens HHD, Crooijmans RPMA, Dibbitts BW, Vereijken A, Okimoto R, Groenen MAM. 2011.
793 Structural variation in the chicken genome identified by paired-end next-generation DNA sequencing of
794 reduced representation libraries. *BMC Genomics* 12.
795 Kim KS, Lee IS, Kim KH, Park J, Kim Y, Choi JH, Choi JS, Jang HJ. 2017. Activation of intestinal
796 olfactory receptor stimulates glucagon-like peptide-1 secretion in enteroendocrine cells and attenuates
797 hyperglycemia in type 2 diabetic mice. *Sci Rep* 7:13978.
798 Korf I. 2004. Gene finding in novel genomes. *BMC Bioinformatics* 5:59.
799 Kotlo K, Anbazhagan AN, Priyamvada S, Jayawardena D, Kumar A, Chen Y, Xia Y, Finn PW, Perkins
800 DL, Dudeja PK, et al. 2020. The olfactory G protein-coupled receptor (Olfir-78/OR51E2) modulates the
801 intestinal response to colitis. *Am J Physiol Cell Physiol* 318:C502-C513.
802 Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. *Nat Methods* 9:357-359.
803 Lawal RA, Al-Atiyat RM, Aljumaah RS, Silva P, Mwacharo JM, Hanotte O. 2018. Whole-Genome
804 Resequencing of Red Junglefowl and Indigenous Village Chicken Reveal New Insights on the Genome
805 Dynamics of the Species. *Front Genet* 9:264.
806 Lawal RA, Martin SH, Vanmechelen K, Vereijken A, Silva P, Al-Atiyat RM, Aljumaah RS, Mwacharo
807 JM, Wu DD, Zhang YP, et al. 2020. The wild species genome ancestry of domestic chickens. *BMC Biol*
808 18:13.
809 Li D, Che T, Chen B, Tian S, Zhou X, Zhang G, Li M, Gaur U, Li Y, Luo M, et al. 2017. Genomic data
810 for 78 chickens from 14 populations. *Gigascience* 6:1-5.
811 Li D, Sun G, Zhang M, Cao Y, Zhang C, Fu Y, Li F, Li G, Jiang R, Han R, et al. 2020. Breeding history
812 and candidate genes responsible for black skin of Xichuan black-bone chicken. *BMC Genomics* 21:511.
813 Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows-Wheeler transform.
814 *Bioinformatics* 25:1754-1760.
815 Li J, Lee M, Davis BW, Lamichhaney S, Dorshorst BJ, Siegel PB, Andersson L. 2020. Mutations
816 Upstream of the TBX5 and PITX1 Transcription Factor Genes Are Associated with Feathered Legs in
817 the Domestic Chicken. *Mol Biol Evol* 37:2477-2486.
818 Li J, Lee MO, Davis BW, Wu P, Hsieh Li SM, Chuong CM, Andersson L. 2021. The crest phenotype in
819 domestic chicken is caused by a 197 bp duplication in the intron of HOXC10. *G3 (Bethesda)* 11.
820 Li M, Chen L, Tian S, Lin Y, Tang Q, Zhou X, Li D, Yeung CKL, Che T, Jin L, et al. 2017. Comprehensive
821 variation discovery and recovery of missing sequence in the pig genome using multiple de novo
822 assemblies. *Genome Res* 27:865-874.

823 Li R, Fu W, Su R, Tian X, Du D, Zhao Y, Zheng Z, Chen Q, Gao S, Cai Y, et al. 2019. Towards the
824 Complete Goat Pan-Genome by Recovering Missing Genomic Segments From the Reference Genome.
825 *Front Genet* 10:1169.

826 Lillie M, Sheng ZY, Honaker CF, Andersson L, Siegel PB, Carlborg O. 2018. Genomic signatures of 60
827 years of bidirectional selection for 8-week body weight in chickens. *Poult Sci* 97:781-790.

828 Ma M, Shen M, Qu L, Dou T, Guo J, Hu Y, Lu J, Li Y, Wang X, Wang K. 2018. Genome-wide association
829 study for carcass traits in spent hens at 72 weeks old. *Italian Journal of Animal Science*:1-6.

830 Malomane DK, Simianer H, Weigend A, Reimer C, Schmitt AO, Weigend S. 2019. The SYNBREED
831 chicken diversity panel: a global resource to assess chicken diversity at high genomic resolution. *BMC*
832 *Genomics* 20:345.

833 McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella K, Altshuler D,
834 Gabriel S, Daly M, et al. 2010. The Genome Analysis Toolkit: a MapReduce framework for analyzing
835 next-generation DNA sequencing data. *Genome Res* 20:1297-1303.

836 Montenegro JD, Golicz AA, Bayer PE, Hurgobin B, Lee H, Chan CK, Visendi P, Lai K, Dolezel J, Batley
837 J, et al. 2017. The pangenome of hexaploid bread wheat. *Plant J* 90:1007-1013.

838 Morejohn GV. 1968a. Breakdown of Isolation Mechanisms in Two Species of Captive Junglefowl
839 (*Gallus Gallus* and *Gallus Sonneratii*). *Evolution* 22:576-582.

840 Morejohn GV. 1968b. Study of Plumage of the Four Species of the Genus *Gallus*. *The Condor* 70:56-65.

841 Pedersen BS, Quinlan AR. 2018. Mosdepth: quick coverage calculation for genomes and exomes.
842 *Bioinformatics* 34:867-868.

843 Pertea M, Kim D, Pertea GM, Leek JT, Salzberg SL. 2016. Transcript-level expression analysis of RNA-
844 seq experiments with HISAT, StringTie and Ballgown. *Nat Protoc* 11:1650-1667.

845 Piegu B, Arensburger P, Beauclair L, Chabault M, Raynaud E, Coustham V, Brard S, Guizard S, Burlot
846 T, Le Bihan-Duval E, et al. 2020. Variations in genome size between wild and domesticated lineages of
847 fowls belonging to the *Gallus gallus* species. *Genomics* 112:1660-1673.

848 Priori D, Colombo M, Clavanzani P, Jansman AJ, Lalles JP, Trevisi P, Bosi P. 2015. The Olfactory
849 Receptor OR51E1 Is Present along the Gastrointestinal Tract of Pigs, Co-Localizes with Enteroendocrine
850 Cells and Is Modulated by Intestinal Microbiota. *PLoS One* 10:e0129501.

851 Qanbari S, Rubin CJ, Maqbool K, Weigend S, Weigend A, Geibel J, Kerje S, Wurmser C, Peterson AT,
852 Brisbin IL, Jr., et al. 2019. Genetics of adaptation in modern chicken. *PLoS Genet* 15:e1007989.

853 Rubin CJ, Zody MC, Eriksson J, Meadows JR, Sherwood E, Webster MT, Jiang L, Ingman M, Sharpe T,
854 Ka S, et al. 2010. Whole-genome resequencing reveals loci under selection during chicken domestication.
855 *Nature* 464:587-591.

856 Schmit MA, Mirakaj V, Stangassinger M, Konig K, Kohler D, Rosenberger P. 2012. Vasodilator
857 Phosphostimulated Protein (VASP) Protects Endothelial Barrier Function During Hypoxia. *Inflammation*
858 35:566-573.

859 Seol D, Ko BJ, Kim B, Chai HH, Lim D, Kim H. 2019. Identification of Copy Number Variation in
860 Domestic Chicken Using Whole-Genome Sequencing Reveals Evidence of Selection in the Genome.
861 *Animals* 9.

862 Sheng Z, Pettersson ME, Hu X, Luo C, Qu H, Shu D, Shen X, Carlborg O, Li N. 2013. Genetic dissection
863 of growth traits in a Chinese indigenous x commercial broiler chicken cross. *BMC Genomics* 14:151.

864 Sherman RM, Forman J, Antonescu V, Puiu D, Daya M, Rafiels N, Boorgula MP, Chavan S, Vergara C,
865 Ortega VE, et al. 2019. Assembly of a pan-genome from deep sequencing of 910 humans of African
866 descent. *Nat Genet* 51:30-35.

867 Singh V, Gowda CP, Singh V, Ganapathy AS, Karamchandani DM, Eshelman MA, Yochum GS, Nighot
868 P, Spiegelman VS. 2020. The mRNA-binding protein IGF2BP1 maintains intestinal barrier function by
869 up-regulating occludin expression. *Journal of Biological Chemistry* 295:8602-8612.

870 Sithu SD, Malovichko MV, Riggs KA, Wickramasinghe NS, Winner MG, Agarwal A, Hamed-Berair RE,
871 Kalani A, Riggs DW, Bhatnagar A, et al. 2017. Atherogenesis and metabolic dysregulation in LDL
872 receptor-knockout rats. *JCI Insight* 2.

873 Smith RW, Cash P, Hogg DW, Buck LT. 2015. Proteomic changes in the brain of the western painted
874 turtle (*Chrysemys picta bellii*) during exposure to anoxia. *Proteomics* 15:1587-1597.

875 Sommer S. 2005. The importance of immune gene variability (MHC) in evolutionary ecology and
876 conservation. *Front Zool* 2:16.

877 Stanke M, Schoffmann O, Morgenstern B, Waack S. 2006. Gene prediction in eukaryotes with a
878 generalized hidden Markov model that uses hints from external sources. *BMC Bioinformatics* 7:62.

879 Stohr N, Kohn M, Lederer M, Glass M, Reinke C, Singer RH, Huttelmaier S. 2012. IGF2BP1 promotes
880 cell migration by regulating MK5 and PTEN signaling. *Genes Dev* 26:176-189.

881 Sun H, Yuan Y, Sun ZL. 2013. Cholesterol Contributes to Diabetic Nephropathy through SCAP-SREBP-
882 2 Pathway. *Int J Endocrinol* 2013:592576.

883 Swinnen G, Goossens A, Pauwels L. 2019. Lessons from Domestication: Targeting Cis-Regulatory
884 Elements for Crop Improvement. *Trends Plant Sci* 24:1065.

885 Tarailo-Graovac M, Chen N. 2009. Using RepeatMasker to identify repetitive elements in genomic
886 sequences. *Curr Protoc Bioinformatics Chapter 4:Unit 4 10*.

887 Tian X, Li R, Fu W, Li Y, Wang X, Li M, Du D, Tang Q, Cai Y, Long Y, et al. 2020. Building a sequence
888 map of the pig pan-genome from multiple de novo assemblies and Hi-C data. *Sci China Life Sci* 63:750-
889 763.

890 Ulfah M, Kawahara-Miki R, Farajallah A, Muladno M, Dorshorst B, Martin A, Kono T. 2016. Genetic
891 features of red and green junglefowls and relationship with Indonesian native chickens Sumatera and
892 Kedu Hitam. *BMC Genomics* 17:320.

893 van der Most PJ, de Jong B, Parmentier HK, Verhulst S. 2011. Trade-off between growth and immune
894 function: a meta-analysis of selection experiments. *Functional Ecology* 25:74-80.

895 Van Laere AS, Nguyen M, Braunschweig M, Nezer C, Collette C, Moreau L, Archibald AL, Haley CS,
896 Buys N, Tally M, et al. 2003. A regulatory mutation in IGF2 causes a major QTL effect on muscle growth
897 in the pig. *Nature* 425:832-836.

898 Varshney RK, Thudi M, Roorkiwal M, He W, Upadhyaya HD, Yang W, Bajaj P, Cubry P, Rathore A, Jian
899 J, et al. 2019. Resequencing of 429 chickpea accessions from 45 countries provides insights into genome
900 diversity, domestication and agronomic traits. *Nat Genet* 51:857-864.

901 Wang M-S, Thakur M, Peng M-S, Jiang Y, Frantz LAF, Li M, Zhang J-J, Wang S, Peters J, Otecko NO,
902 et al. 2020a. 863 genomes reveal the origin and domestication of chicken. *Cell Research* 30:693-701.

903 Wang MS, Li Y, Peng MS, Zhong L, Wang ZJ, Li QY, Tu XL, Dong Y, Zhu CL, Wang L, et al. 2015.
904 Genomic Analyses Reveal Potential Independent Adaptation to High Altitude in Tibetan Chickens. *Mol
905 Biol Evol* 32:1880-1889.

906 Wang MS, Thakur M, Peng MS, Jiang Y, Frantz LAF, Li M, Zhang JJ, Wang S, Peters J, Otecko NO, et
907 al. 2020b. 863 genomes reveal the origin and domestication of chicken. *Cell Res* 30:693-701.

908 Wang Y, Bu L, Cao X, Qu H, Zhang C, Ren J, Huang Z, Zhao Y, Luo C, Hu X, et al. 2020. Genetic
909 Dissection of Growth Traits in a Unique Chicken Advanced Intercross Line. *Front Genet* 11:894.

910 Wang Y, Cao X, Luo C, Sheng Z, Zhang C, Bian C, Feng C, Li J, Gao F, Zhao Y, et al. 2020. Multiple

- 911 ancestral haplotypes harboring regulatory mutations cumulatively contribute to a QTL affecting chicken
 912 growth traits. *Commun Biol* 3:472.
- 913 Wang Z, Qu L, Yao J, Yang X, Li G, Zhang Y, Li J, Wang X, Bai J, Xu G, et al. 2013. An EAV-HP
 914 insertion in 5' Flanking region of *SLCO1B3* causes blue eggshell in the chicken. *PLoS Genet* 9:e1003183.
- 915 Warner C, Meeker D, Rothschild M. (2.092 co-authors). 1987. Genetic control of immune responsiveness:
 916 a review of its use as a tool for selection for disease resistance. *Journal of animal science* 64:394-406.
- 917 Whitlock MC. 2000. Fixation of new alleles and the extinction of small populations: drift load, beneficial
 918 alleles, and sexual selection. *Evolution* 54:1855-1861.
- 919 Wood DE, Lu J, Langmead B. 2019. Improved metagenomic analysis with Kraken 2. *Genome Biol*
 920 20:257.
- 921 Wright D, Boije H, Meadows JR, Bed'hom B, Gourichon D, Vieaud A, Tixier-Boichard M, Rubin CJ,
 922 Imsland F, Hallbook F, et al. 2009. Copy number variation in intron 1 of *SOX5* causes the Pea-comb
 923 phenotype in chickens. *PLoS Genet* 5:e1000512.
- 924 Yang J, Benyamin B, McEvoy BP, Gordon S, Henders AK, Nyholt DR, Madden PA, Heath AC, Martin
 925 NG, Montgomery GW, et al. 2010. Common SNPs explain a large proportion of the heritability for
 926 human height. *Nat Genet* 42:565-569.
- 927 Yao W, Li G, Zhao H, Wang G, Lian X, Xie W. 2015. Exploring the rice dispensable genome using a
 928 metagenome-like assembly strategy. *Genome Biol* 16:187.
- 929 Zhang L, Wan YC, Zhang ZH, Jiang Y, Gu ZY, Ma XL, Nie SP, Yang J, Lang JH, Cheng WJ, et al. 2021.
 930 *IGF2BP1* overexpression stabilizes *PEG10* mRNA in an m6A-dependent manner and promotes
 931 endometrial cancer progression. *Theranostics* 11:1100-1114.
- 932 Zhang Y, Ma KL, Ruan XZ, Liu BC. 2016. Dysregulation of the Low-Density Lipoprotein Receptor
 933 Pathway Is Involved in Lipid Disorder-Mediated Organ Injury. *Int J Biol Sci* 12:569-579.
- 934 Zhang Y, Wang Y, Li Y, Wu J, Wang X, Bian C, Tian Y, Sun G, Han R, Liu X, et al. 2020. Genome-wide
 935 association study reveals the genetic determinism of growth traits in a Gushi-Anka F2 chicken population.
 936 *Heredity (Edinb)*.
- 937 Zhao PJ, Li JH, Kang HM, Wang HF, Fan ZY, Yin ZJ, Wang JF, Zhang Q, Wang ZQ, Liu JF. 2016.
 938 Structural Variant Detection by Large-scale Sequencing Reveals New Evolutionary Evidence on Breed
 939 Divergence between Chinese and European Pigs. *Scientific Reports* 6.
- 940 Zhao Q, Feng Q, Lu H, Li Y, Wang A, Tian Q, Zhan Q, Lu Y, Zhang L, Huang T, et al. 2018. Pan-genome
 941 analysis highlights the extent of genomic variation in cultivated and wild rice. *Nat Genet* 50:278-284.
- 942 Zhou Z, Li M, Cheng H, Fan W, Yuan Z, Gao Q, Xu Y, Guo Z, Zhang Y, Hu J, et al. 2018. An intercross
 943 population study reveals genes associated with body size and plumage color in ducks. *Nat Commun*
 944 9:2648.
- 945 Zhu S, Wang JZ, Chen, He YT, Meng N, Chen M, Lu RX, Chen XH, Zhang XL, Yan GR. 2020. An
 946 oncopeptide regulates m(6)A recognition by the m(6)A reader *IGF2BP1* and tumorigenesis. *Nat*
 947 *Commun* 11:1685.
- 948 Zillikens MC, Demissie S, Hsu YH, Yerges-Armstrong LM, Chou WC, Stolk L, Livshits G, Broer L,
 949 Johnson T, Koller DL, et al. 2017. Large meta-analysis of genome-wide association studies identifies
 950 five loci for lean body mass. *Nat Commun* 8:80.
- 951 Zimin AV, Marcais G, Puiu D, Roberts M, Salzberg SL, Yorke JA. 2013. The MaSuRCA genome
 952 assembler. *Bioinformatics* 29:2669-2677.

953 **Figure legends**

954 **Figure 1. Pan-genome of chicken. a** Geographical distribution of samples used for
 955 pan-genome construction. **b** Pan-genome gene classification. **c** Word cloud of the Gene
 956 Ontology (GO) enrichment of biological process for variable genes. **d** Pan-genome
 957 modelling. The pan-genome modelling shows no more dramatic increases when the
 958 number of accession genomes is over 220, indicating that selected individuals were
 959 sufficient to capture the majority of PAVs within *Gallus gallus*. Upper and lower lines
 960 represent the pan-genome number and core-genome number, respectively.

961 **Figure 2. Distribution of gene PAV. a** The heatmap shows the PAV of variable genes
 962 within wild relatives, native breeds and commercial breeds. **b** The principal component
 963 analysis of chicken breeds based on gene PAV. Wild: wild relatives (red jungle fowls);
 964 Native, native breeds; commercial breeds consist of two broiler breeds (BRA and BRB)
 965 and two layer breeds (BL and WL). **c** Neighbor-Joining phylogenetic tree constructed
 966 based on gene PAV matrix.

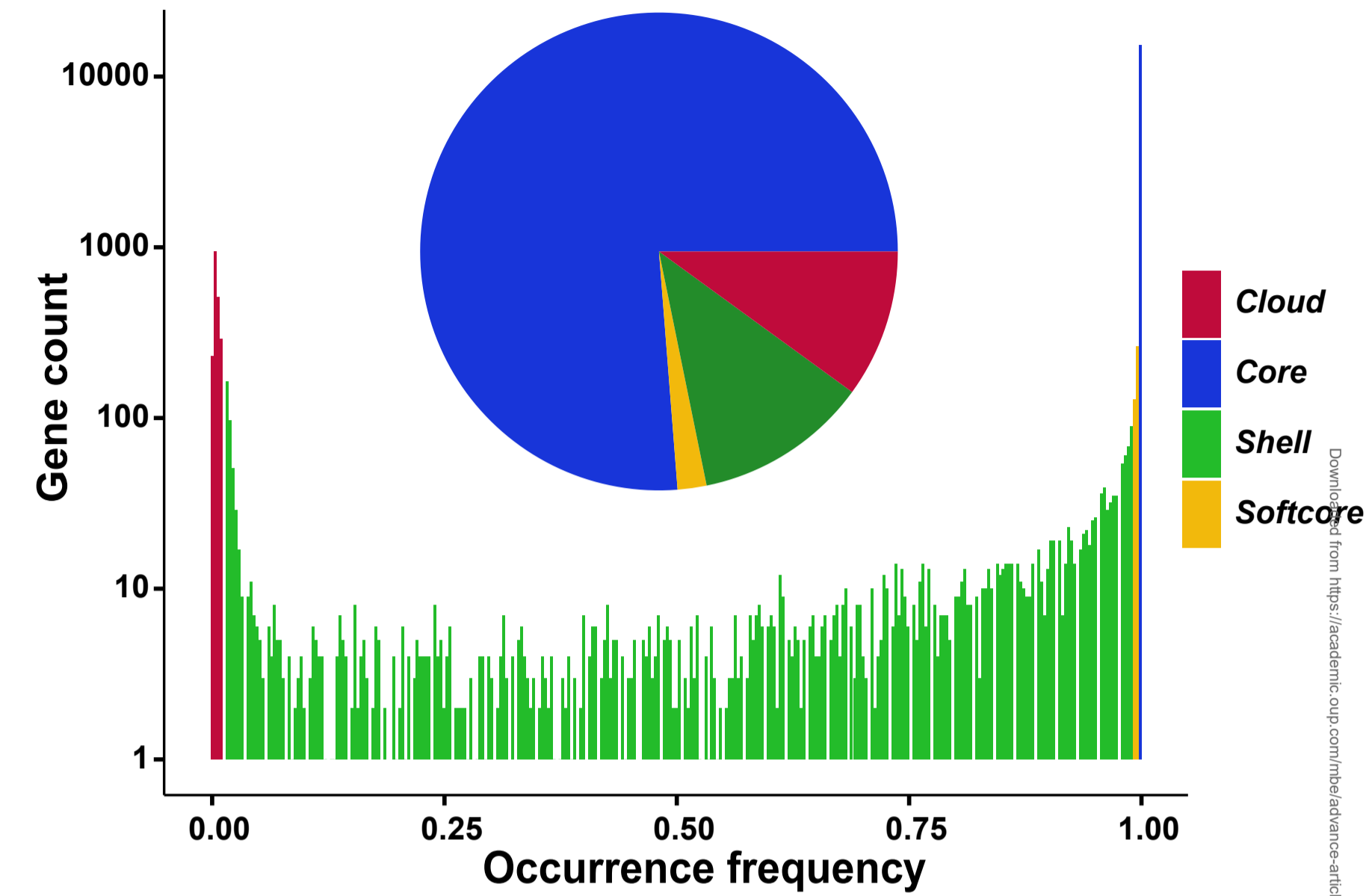
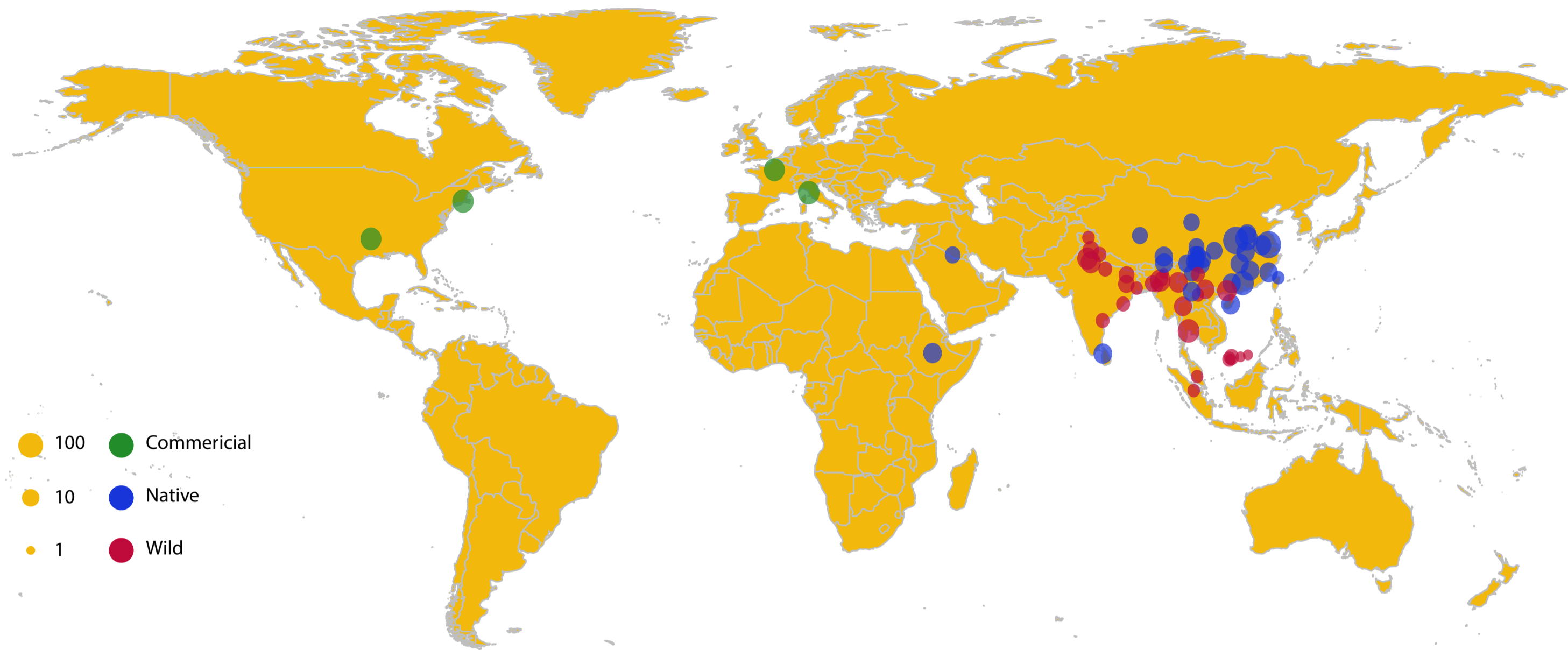
967 **Figure 3. Change of PAV frequency in promoter region during breeding and PAV-**
 968 **based GWAS. a, b, c,** Scatter plots showing gene occurrence frequencies in Native
 969 breeds and Com (commercial) breeds for 0-1 kb (**a**), 1-2 kb (**b**) and 2-3 kb (**c**) upstream
 970 promoter regions, respectively. **d, e, f,** Manhattan plots showing significant promoter
 971 region PAVs associated with 151 traits for 0-1 kb (**d**), 1-2 kb (**e**), and 2-3 kb (**f**) upstream
 972 promoter regions. All association analysis result was plotted according to the physical
 973 location and p-value, with each dot representing an association analysis result. The
 974 upper and lower dashed line represents the significant and suggestive thresholds,
 975 respectively. CW1, claw weight; CR, the ratio of claw weight to body weight; DPW,
 976 double pinion weight; SEW, semi-evisceration weight.

977 **Figure 4. Structure and frequency of the three alleles in *IGF2BP1* promoter region.**
 978 **a** Genomic structure of three alleles in *IGF2BP1* promoter region in relation to
 979 evolutionarily conserved elements (77 vertebrates basewise PhyloP conservation score).
 980 Variant alleles in the promoter region of *IGF2BP1* include wild type (W) and two
 981 mutant alleles (L1 and L2). The conserved elements are indicated by red arrows. *Asp-*
 982 *F*, *2k-F* and *Asp-R* are the PCR primers for the identification of the allelic type. **b** Allelic

983 frequency of *IGF2BP1* promoter region in the validated population by allelic-specific
984 PCR genotyping. PCR product sizes of W, L1 and L2 are 2345 bp, 290 bp and 791 bp,
985 respectively. The gel shows the six genotypes derived from the combinations of the
986 three alleles.

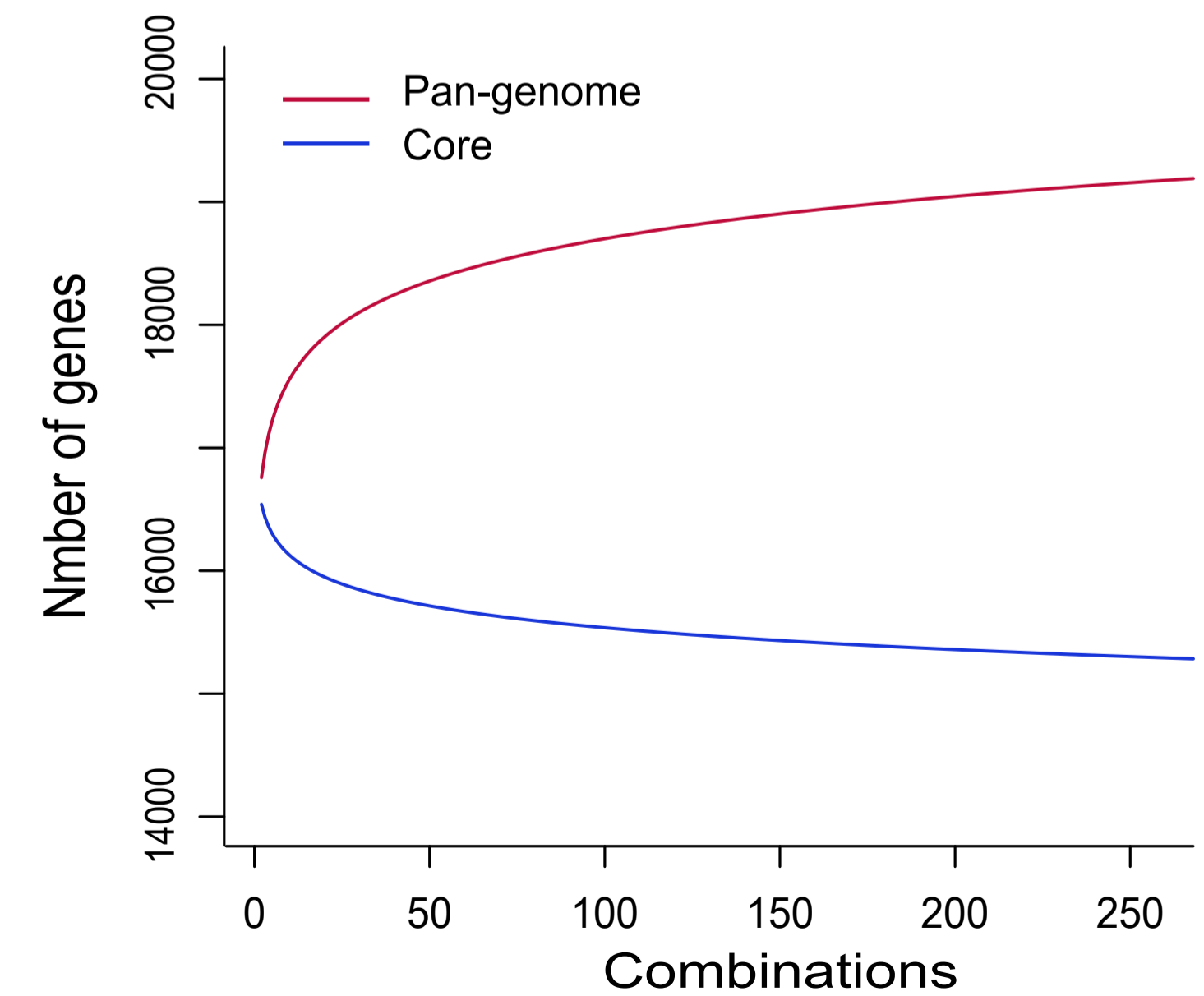
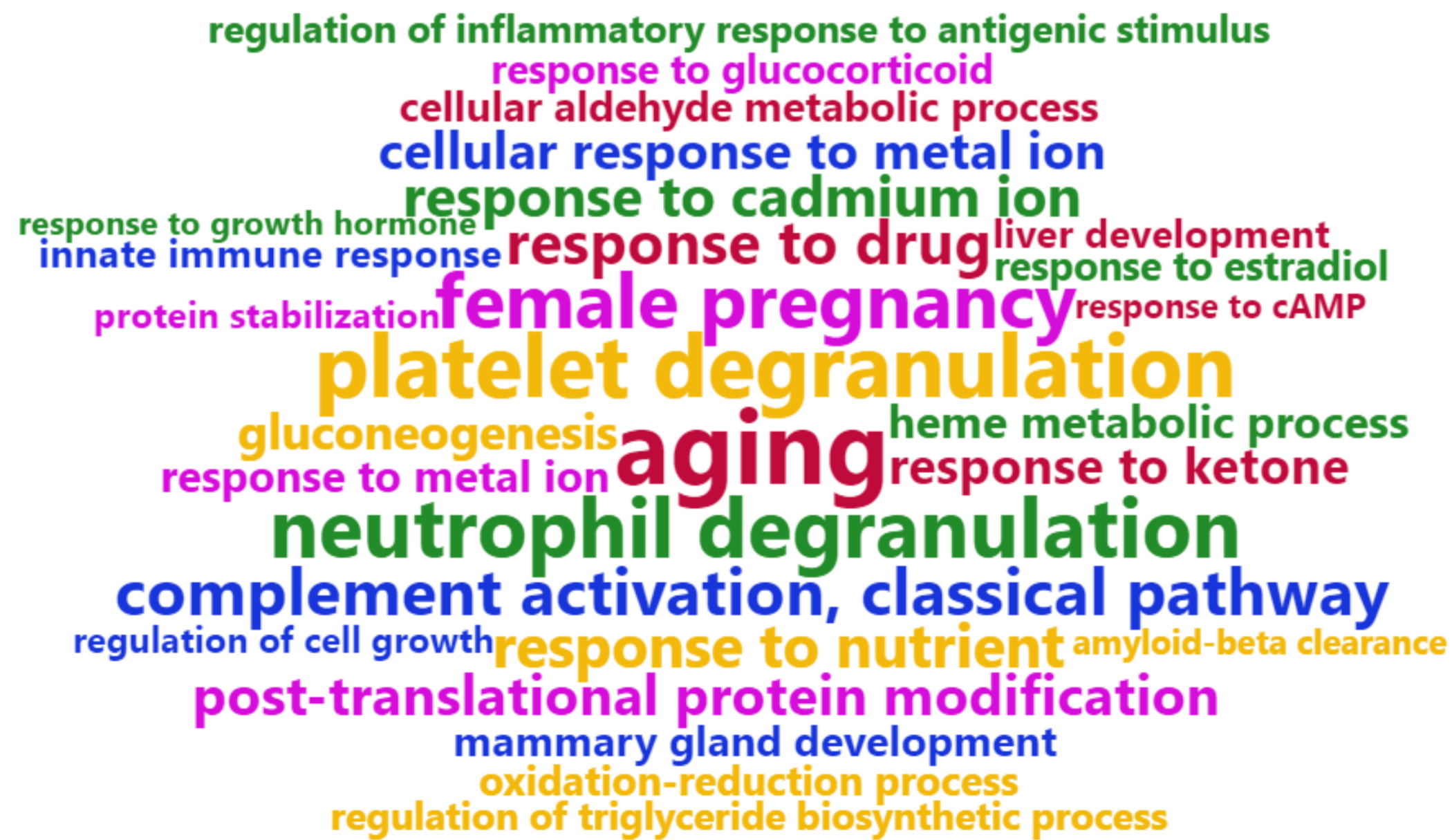
987 **Figure 5. Single-marker genotype association of *IGF2BP1* promoter region in the**
988 **validated Gushi×Anak F2 population with 734 individuals.** Eight representative
989 association events were included and others were showed in Supplementary Figure S10.
990 The number in the bracket is the proportion of phenotype variance explained by
991 *IGF2BP1* loci. CW1, claw weight; CR, the ratio of claw weight to body weight. SL12,
992 shank length; BBL12, breast bone length; DPW, double pinion weight; SEW, semi-
993 evisceration weight; CW, carcass weight; LW, leg weight. All traits were phenotyped at
994 12 weeks of age.

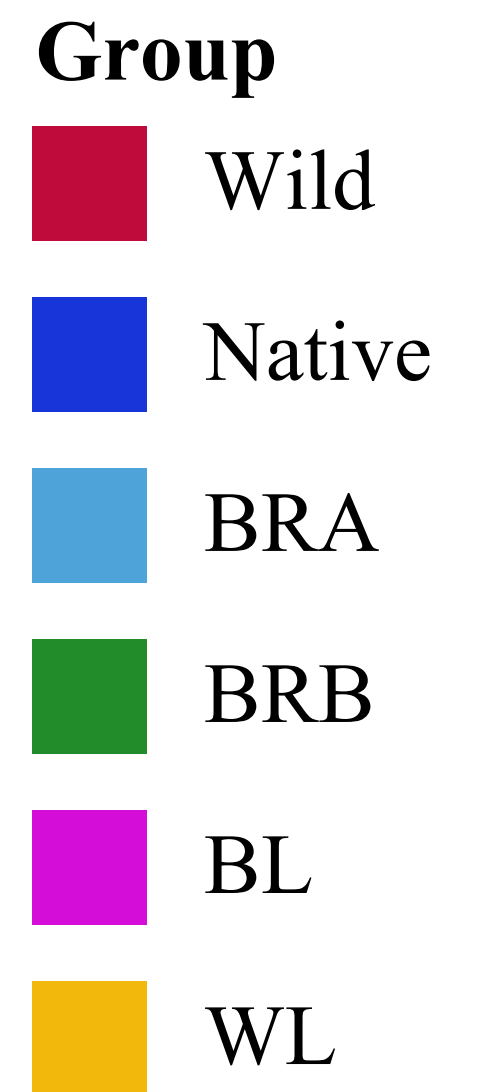
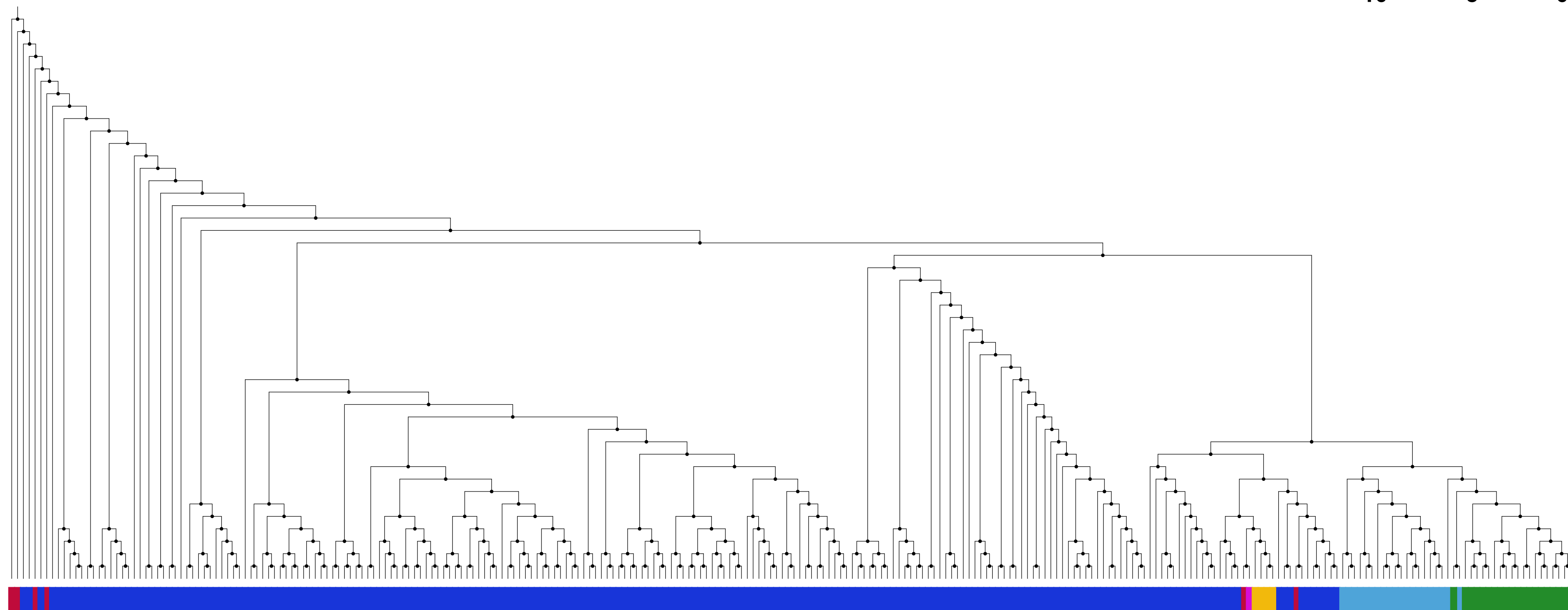
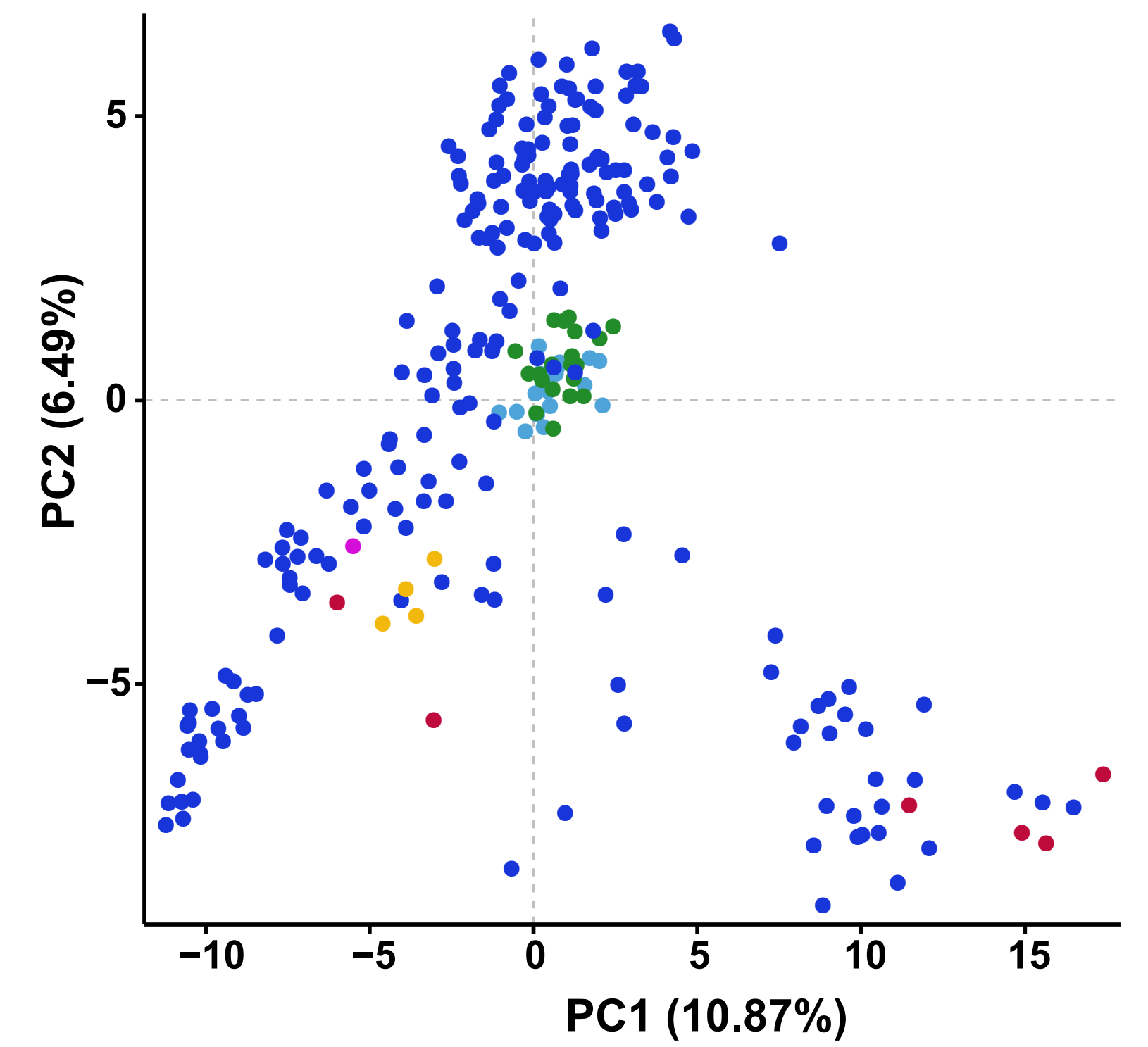
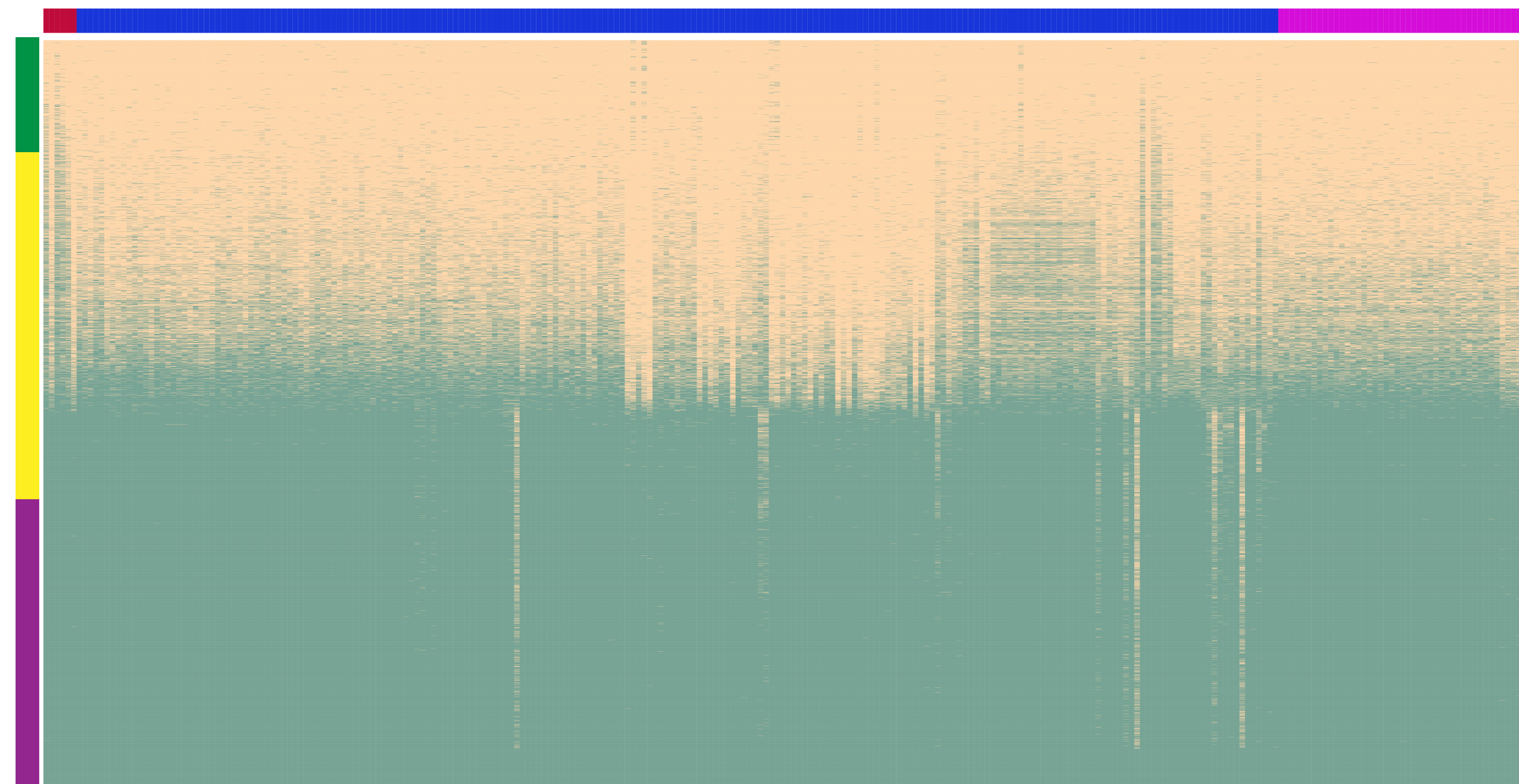
995 **Figure 6. Comparison of transcriptional activity and expression among three**
996 ***IGF2BP1* genotypes.** **a** Comparison of transcriptional activity among different
997 *IGF2BP1* promoter region in chicken DF-1 cells. Left shows the constructions of the
998 inserted fragment into the pGL3-Basic plasmid. Significance of two-tailed Student's *t*-
999 test: **, $p < 0.01$; ***, $p < 0.001$. **b** Comparison of mRNA expression of *IGF2BP1*
1000 between L1L1 (Ross 308) and WW (Gushi) chickens in five tissues at 6 weeks of age.
1001 Breast, breast muscle; Leg, leg muscle. P-values were calculated using a two-tailed
1002 Student's *t*-test. **c** Comparison of mRNA expression of *IGF2BP1* between L1L1, L2L2
1003 and WW in an *IGF2BP1* genotype segregating population at 3 weeks of age.
1004

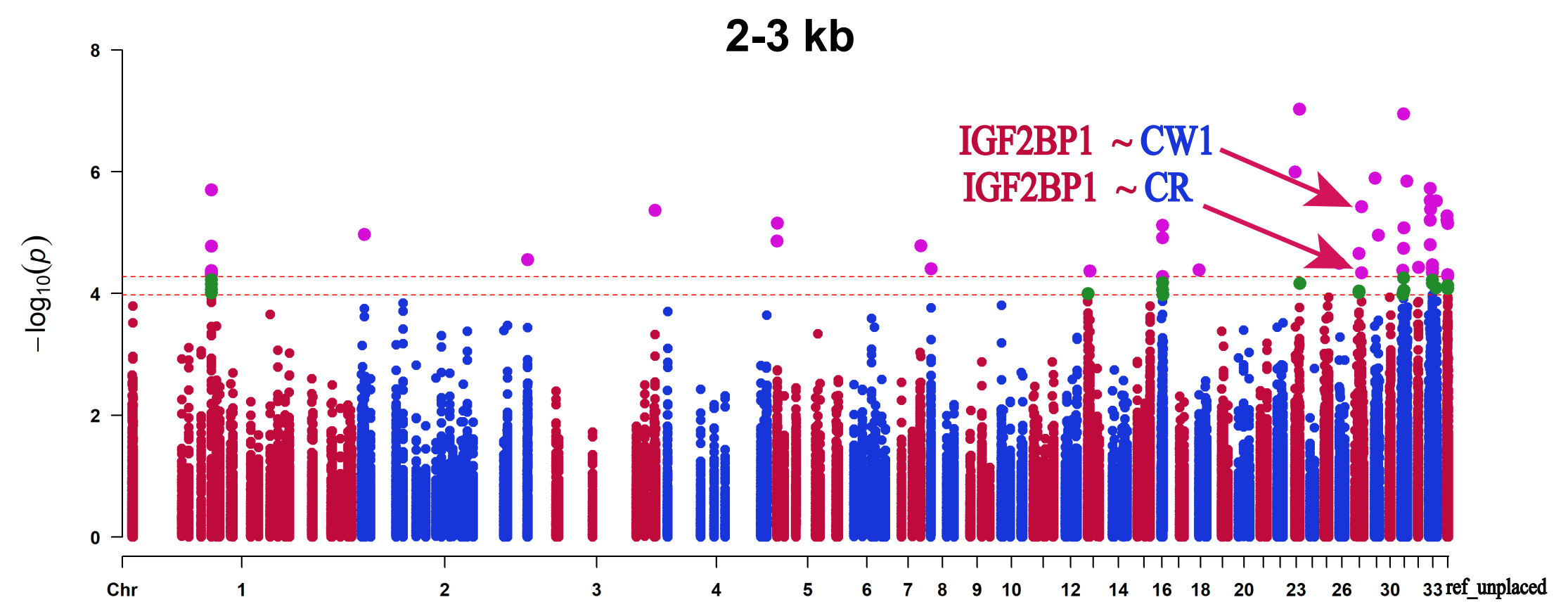
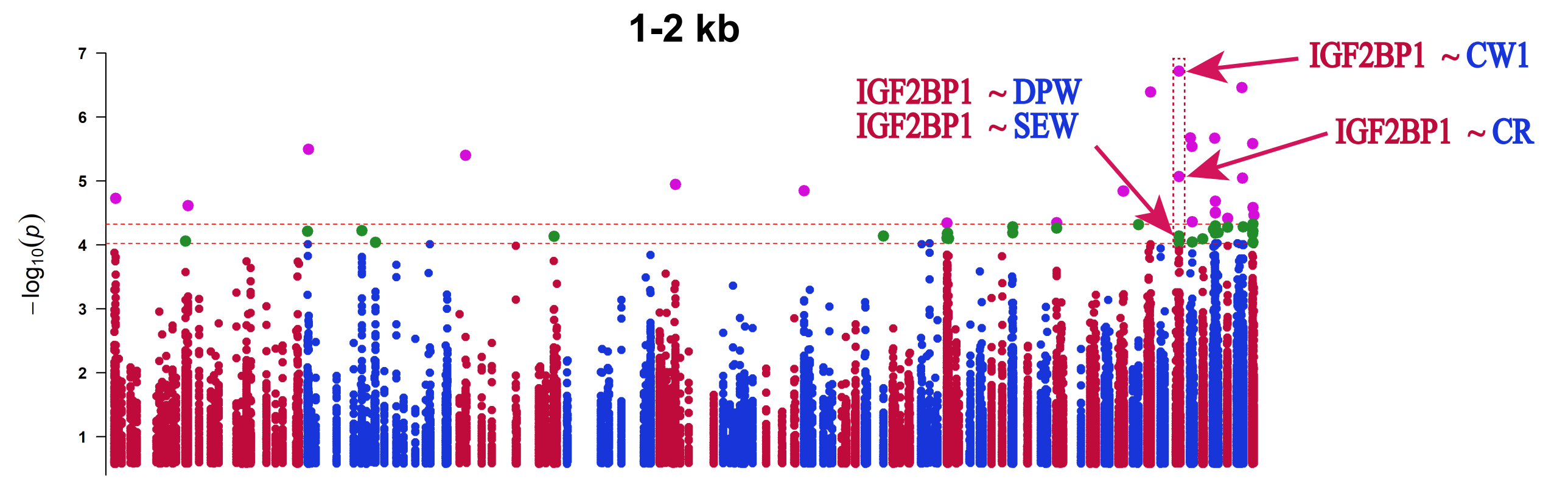
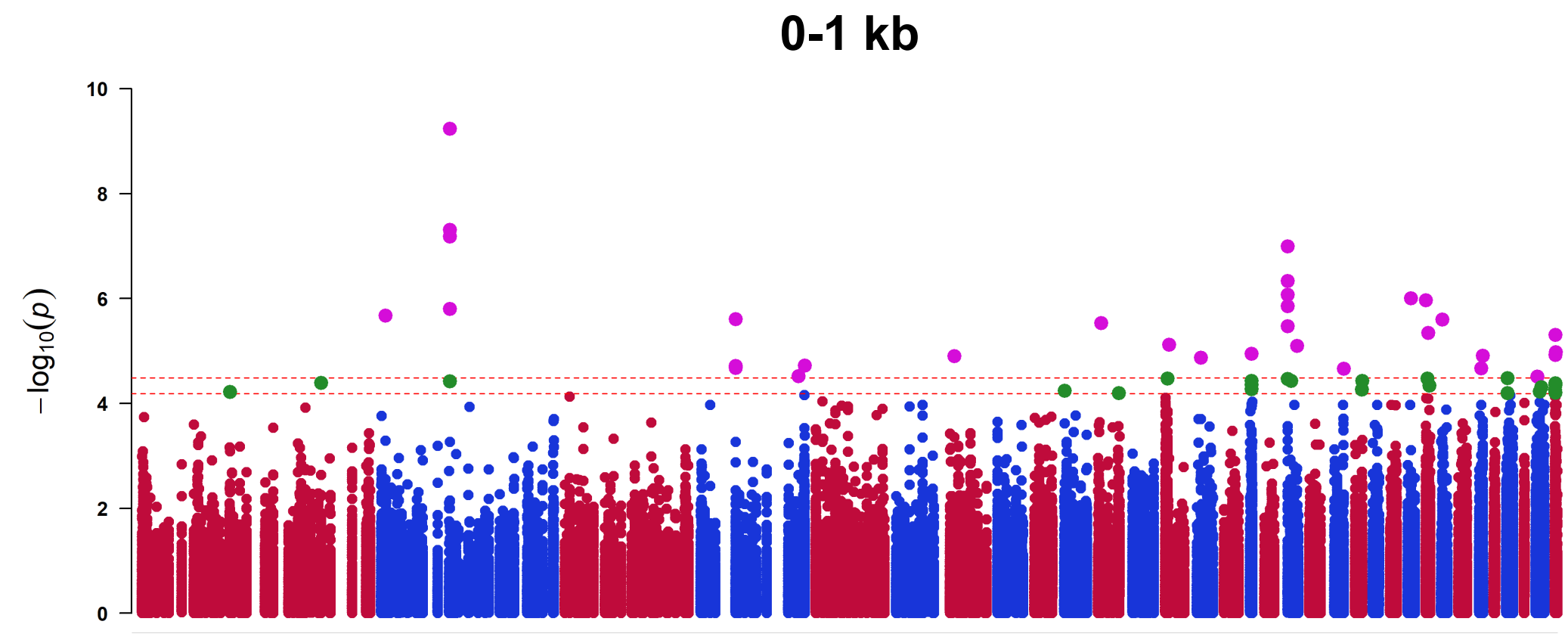
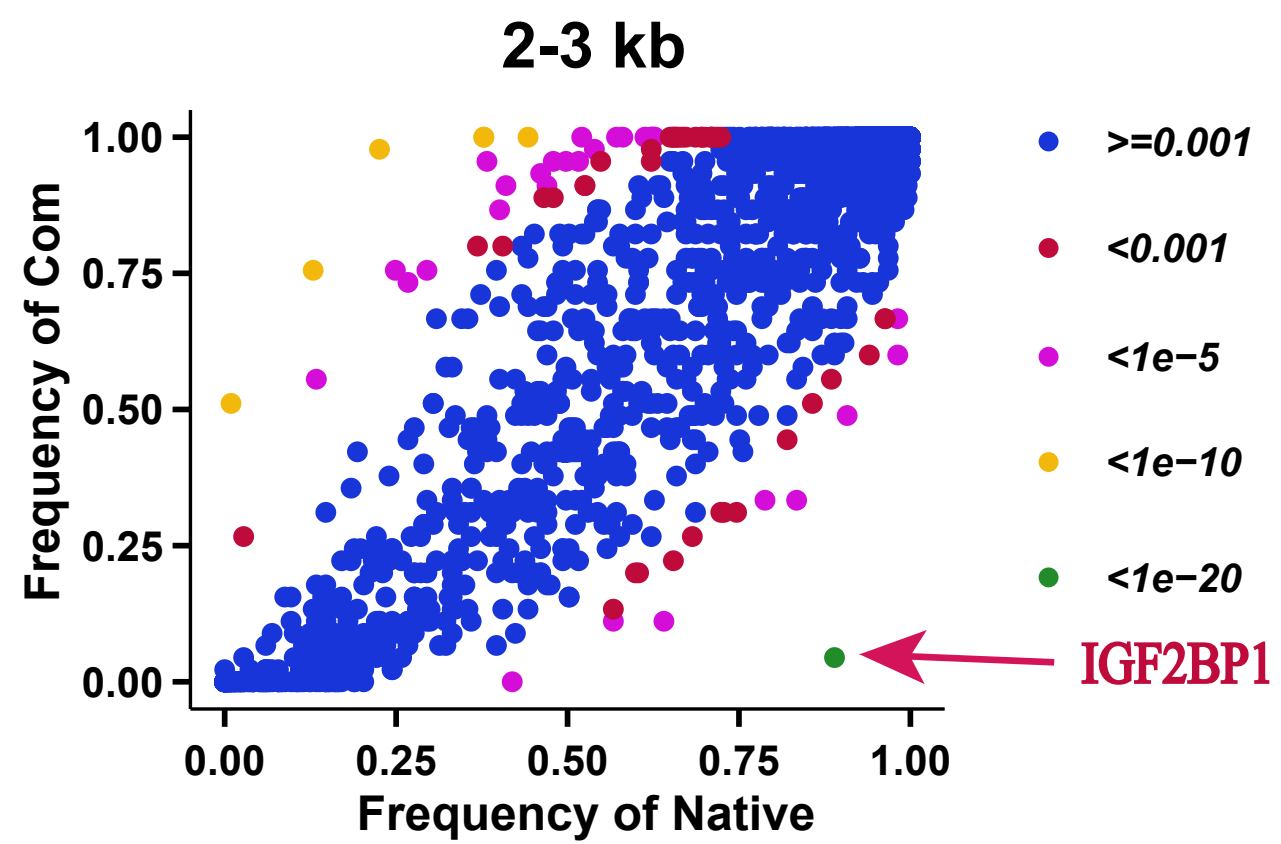
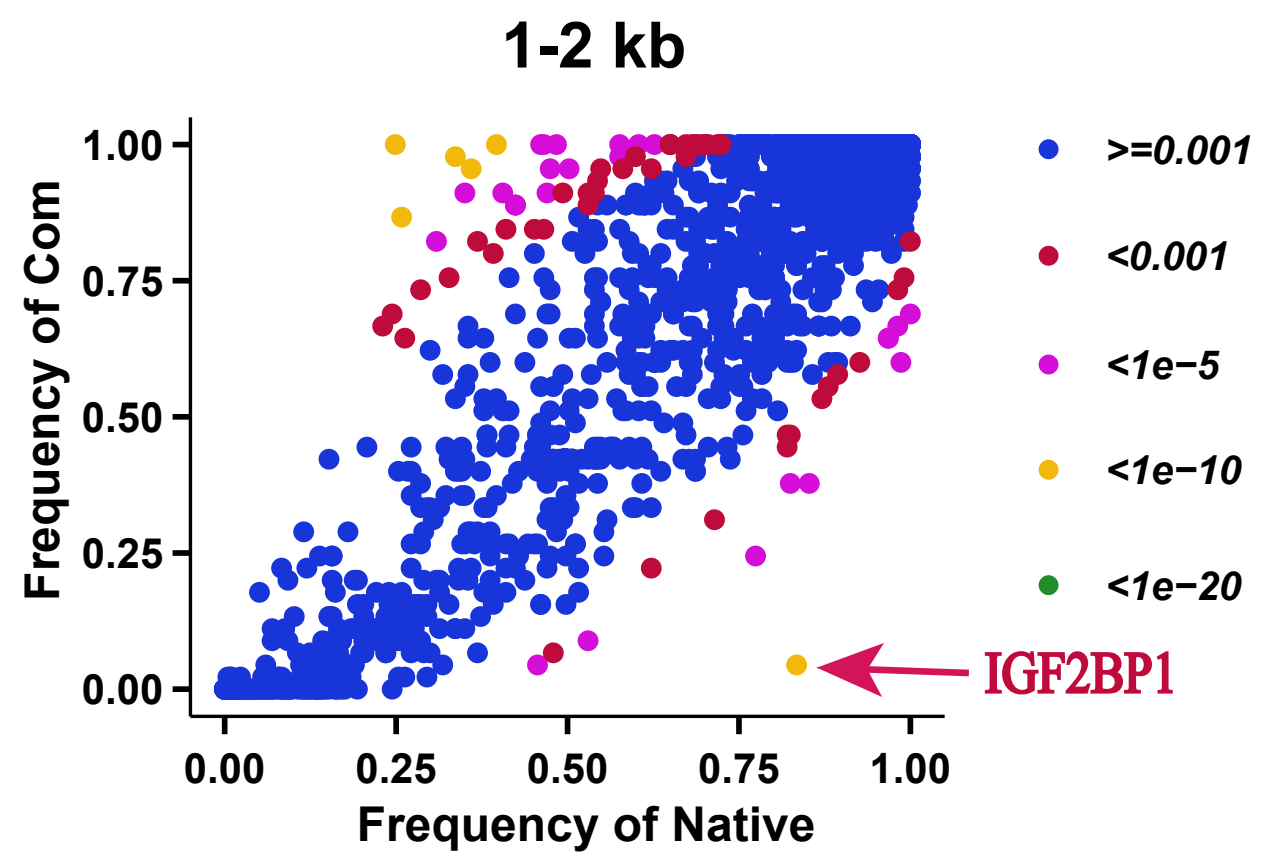
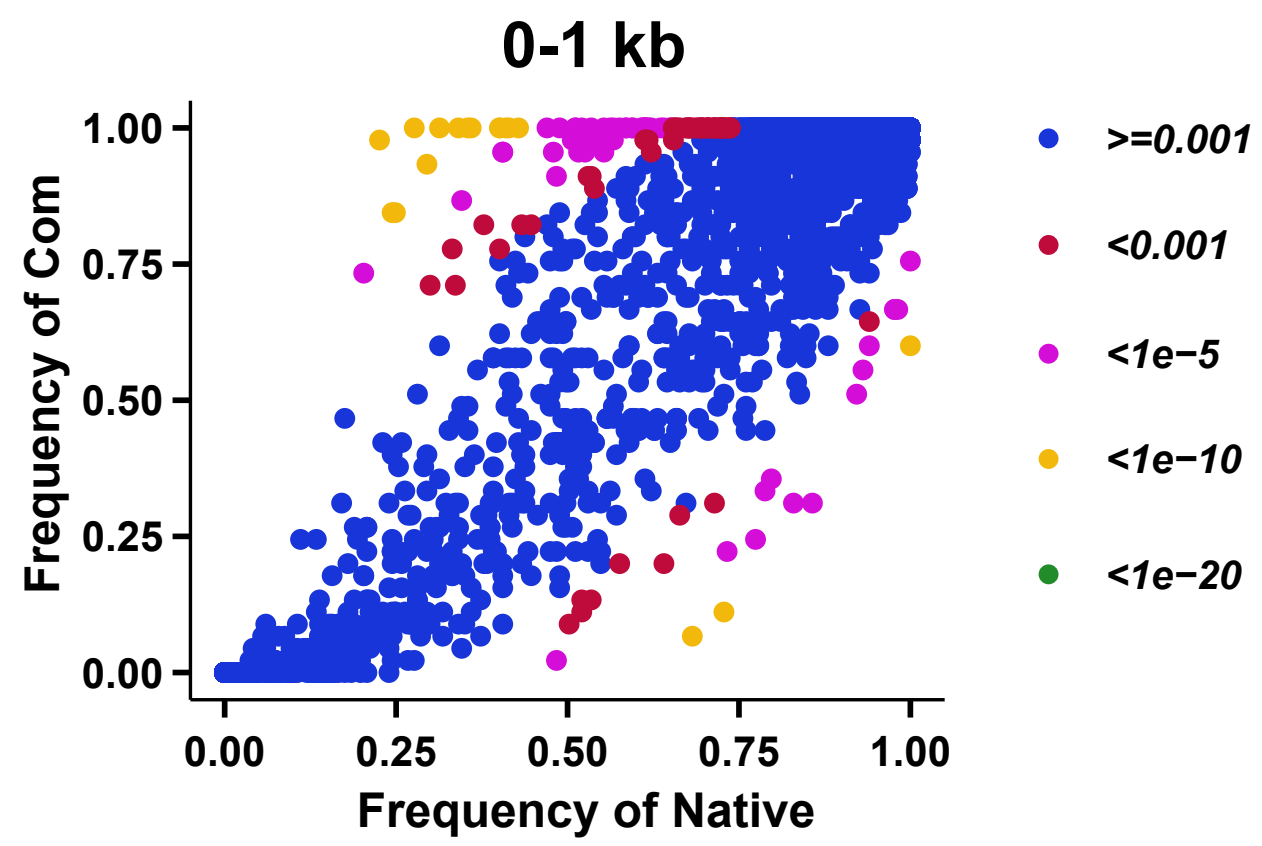


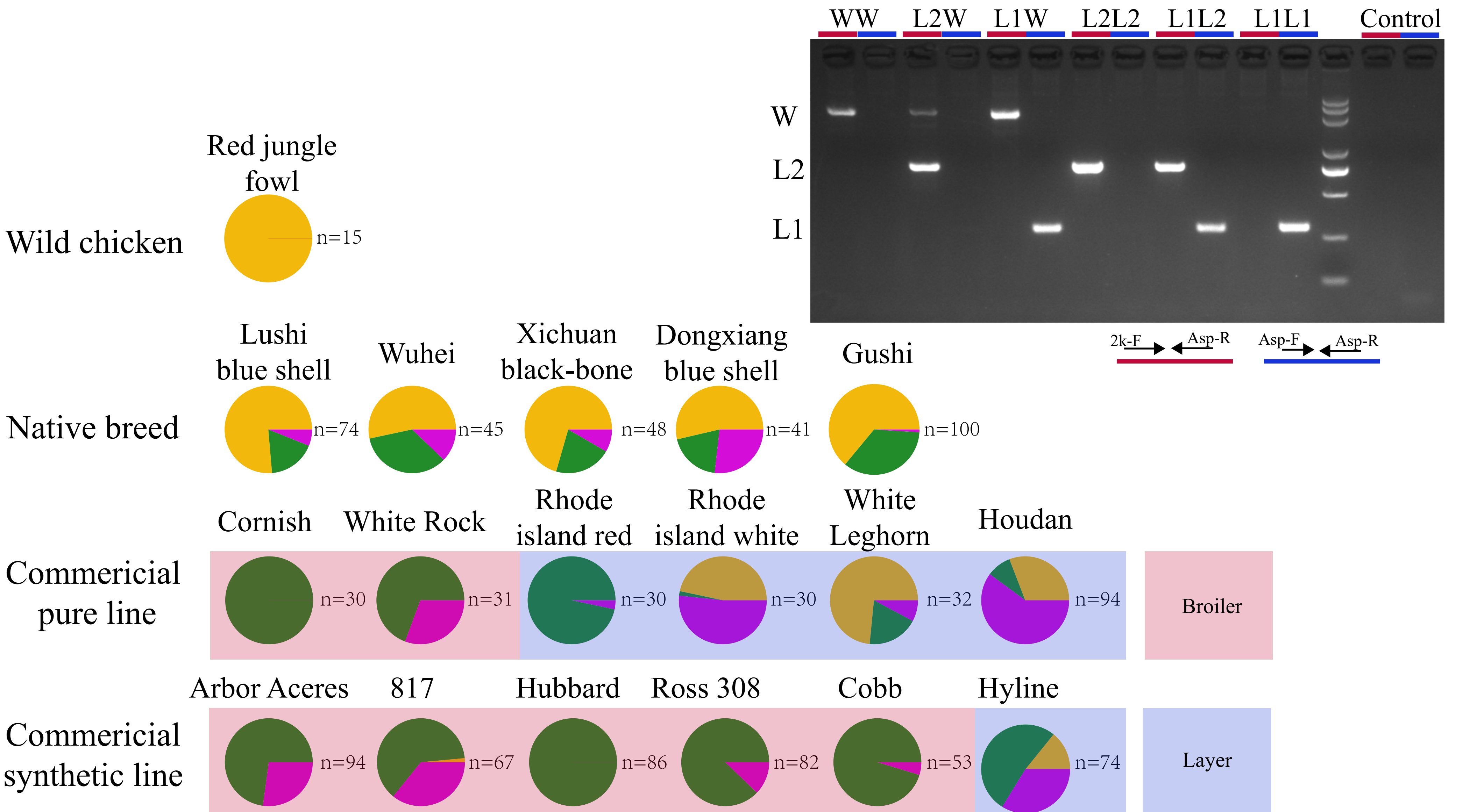
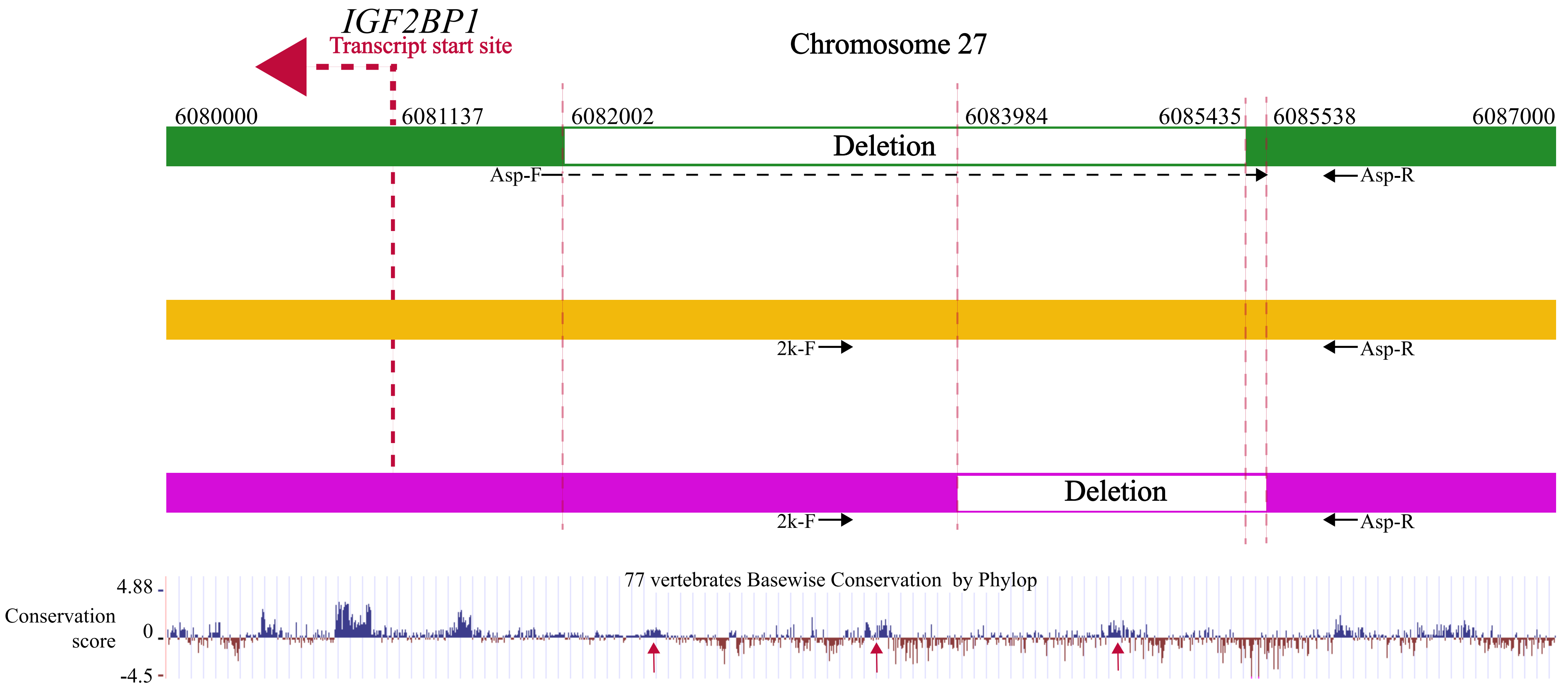
	Core	Softcore	Shell	Cloud
Ref	15042	300	518	0
Non-ref	163	91	1833	1976
Total	15205	391	2351	1976

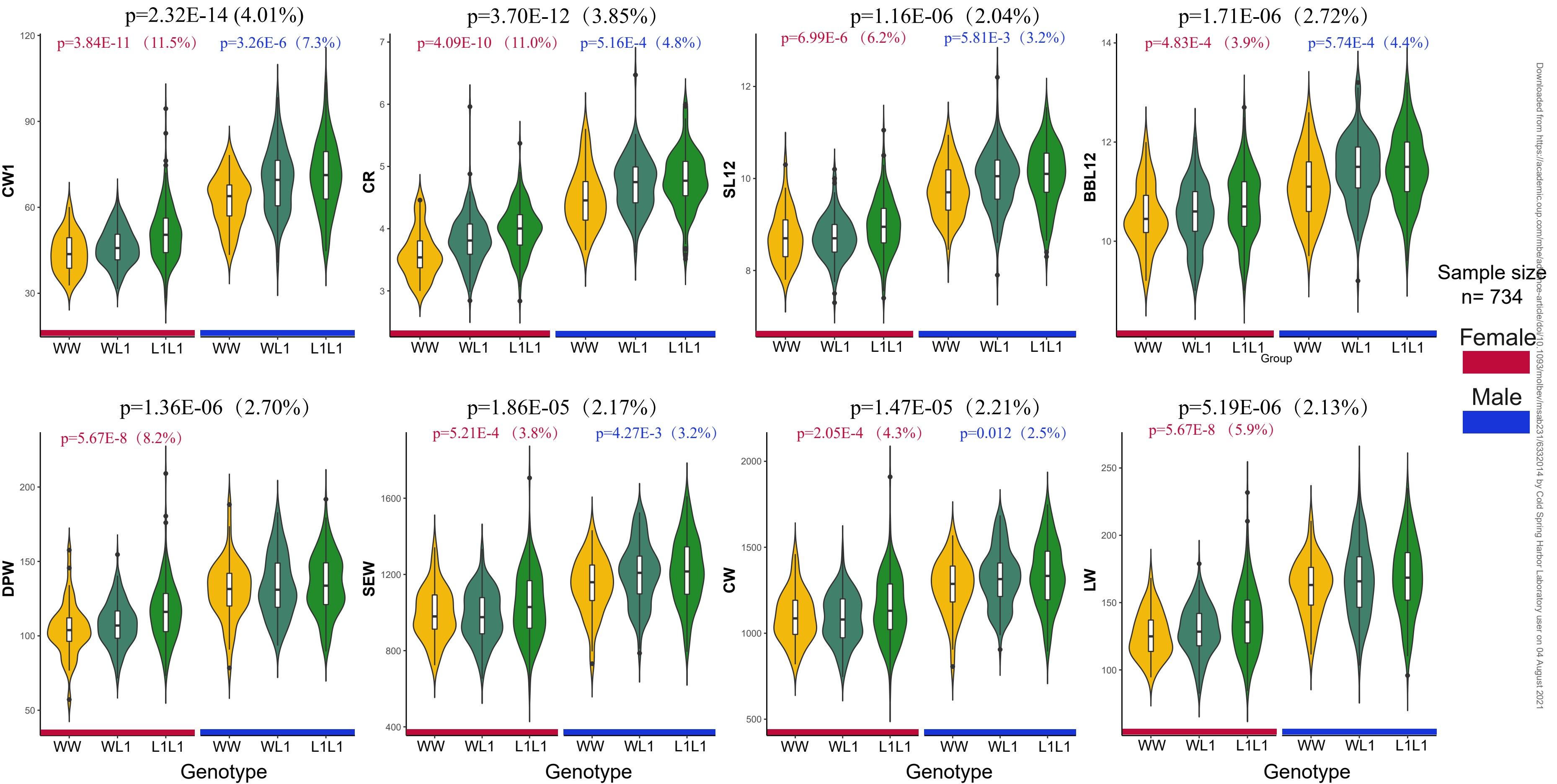
Downloaded from https://academic.oup.com/iad/advance-article/doi/10.1093/iad/iaab231/6322014 by Cold Spring Harbor Laboratory user on 04 August 2021

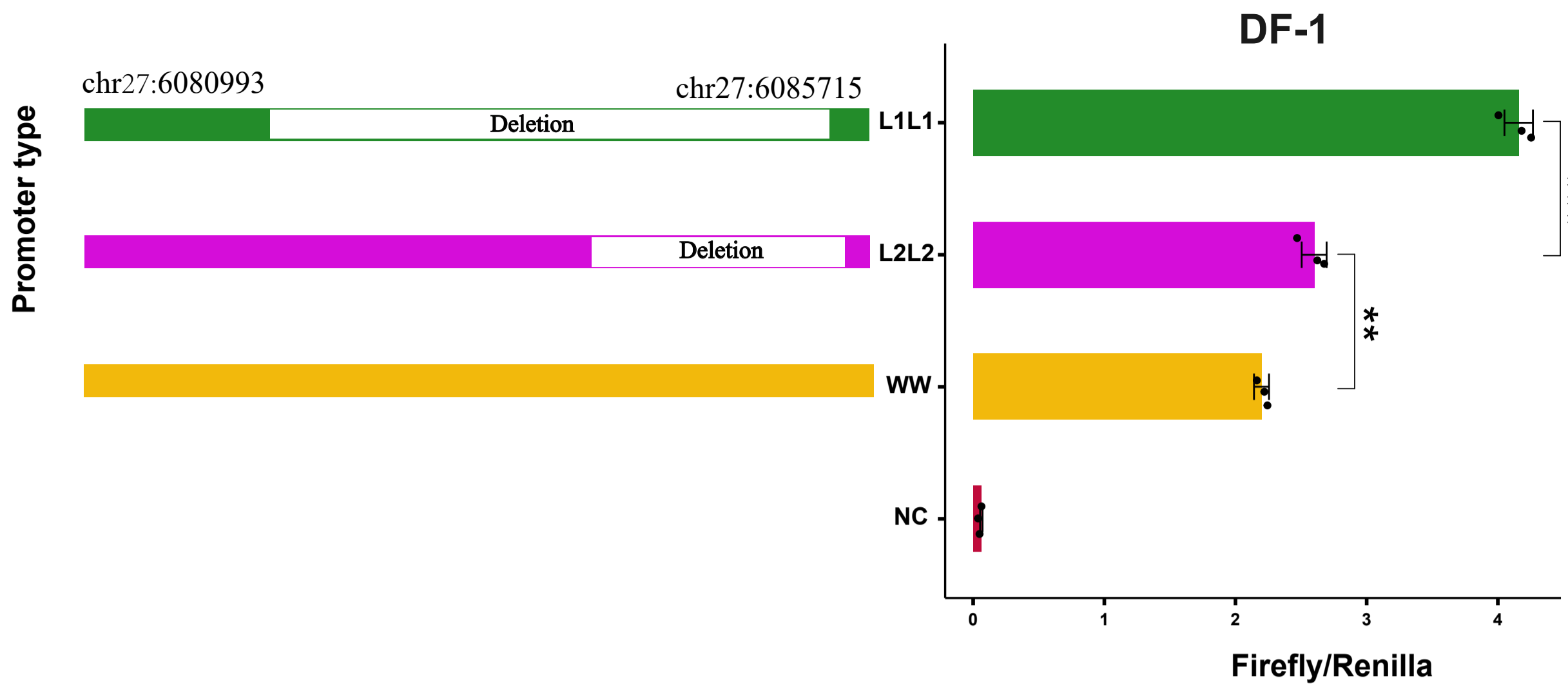
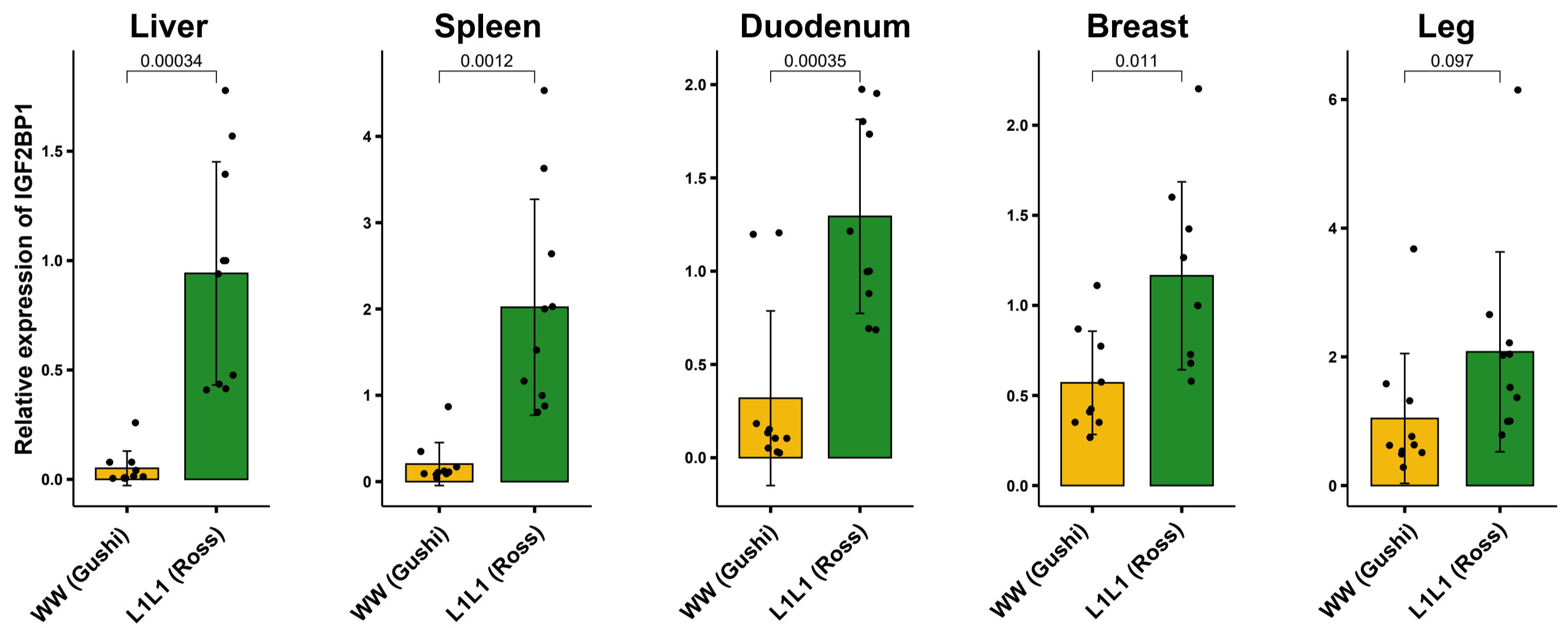










a**b****c**