

# 1 **Kinship analysis and pedigree reconstruction by RAD sequencing in cattle**

2 Yiming Xu<sup>1,\*</sup>, Wanqiu Wang<sup>2,\*</sup>, Minjie Xu<sup>3,\*</sup>, Binhu Wang<sup>2</sup>, Jiefeng Huang<sup>4</sup>, Yingsong Wu<sup>4</sup>,

3 Yongzhong Xie<sup>1</sup>, Jianbo Jian<sup>2,5,6#</sup>

4 <sup>1</sup> Animal Husbandry and Aquatic Affairs Center, Lianyuan City, Hunan Province, China.

5 <sup>2</sup> BGI Genomics, Shenzhen, China.

6 <sup>3</sup> People's Government of Shexian County, Hebei Province, China.

7 <sup>4</sup> Loudi Municipal Bureau of Agriculture and Rural Affairs, Loudi City, Hunan Province,

8 China.

9 <sup>5</sup>Department of Biotechnology and Biomedicine, Technical University of Denmark, Lyngby,  
10 Denmark.

11 <sup>6</sup>Marine Biology Institute, Shantou University, Shantou, China.

12

13 #Correspondence: Jianbo Jian (jbjian@126.com);

14 \*These authors contributed equally: Yiming Xu, Wanqiu Wang (ORCID: 0000-0002-2748-

15 5867), Minjie Xu.

16

17 **Abstract:** Kinship and pedigree information, used for estimating inbreeding, heritability,

18 selection, and gene flow, is useful for breeding and animal conservation. However, as the size

19 of the crossbred population increases, inaccurate generation and parentage recoding in livestock

20 farms increases. Restriction-site-associated DNA sequencing (RAD-Seq) is a cost-effective

21 platform for single nucleotide polymorphism (SNP) discovery and genotyping. Here, we

22 performed a kinship analysis and pedigree reconstruction for Angus and Xiangxi yellow cattle,

23 which benefit from good meat quality and yields, providing a basis for livestock management.

24 A total of 975 cattle, including 923 offspring with 24 known sires and 28 known dams, were

25 sampled and subjected to SNP discovery and genotyping. The identified SNPs panel included

26 7305 SNPs capturing the maximum difference between paternal and maternal genome  
27 information allowing us to distinguish between the F1 and F2 generation with 90% accuracy.  
28 In addition, parentage assignment software based on different strategies verified that the cross-  
29 assignments. In conclusion, we provided a low-cost and efficient SNP panel for kinship  
30 analyses and the improvement of local genetic resources, which are valuable for breed  
31 improvement, local resource utilization, and conservation.

32 **Keywords:** crossbreeding, RAD-seq, generation classification, parentage assignment

## 33 **Introduction**

34 The selection of superior individuals as parents in crossbreeding systems is a general tool for  
35 genetic improvement (Gregory et al. 1980). Genetic selection is a new approach to selective  
36 breeding based on high-density markers covering the entire genome; it shortens the breeding  
37 cycle and reduces breeding costs. Predicting progeny from parental lines and selecting the best  
38 crosses rely on the accuracy of the pedigree; however, parentage assignment is expensive and  
39 laborious in normal livestock production systems. Therefore, the development of reliable  
40 genetic markers for kinship analyses and pedigree reconstruction is necessary to meet industry  
41 demand and improve livestock management systems.

42  
43 SNPs are abundant, widespread in the genome, and follow simple models of evolution (Morin  
44 et al. 2004). There are various genotyping approaches, including bovine SNP chip, restriction-  
45 site-associated DNA sequencing (RAD-seq) (Baird et al. 2008), multiplexed shotgun  
46 genotyping (Andolfatto et al. 2011), exome sequencing (Ng et al. 2009), and whole-genome  
47 resequencing (WGS). Traditional SNP genotyping methods, which necessitate specific probes  
48 or primers tailored to each SNP of interest, can incur substantial costs, particularly for large-  
49 scale genotyping efforts. Moreover, these processes are often labor-intensive and time-

50 consuming, especially in the absence of high-throughput methodologies. Additionally, they  
51 may necessitate sophisticated bioinformatics for data analysis, further augmenting the  
52 associated costs. In contrast, RAD sequencing generally employs a restricted set of restriction  
53 enzymes and barcodes, which can significantly reduce reagent costs compared to individual  
54 SNP assays. Despite the complexity of next-generation sequencing (NGS) data, a plethora of  
55 tools and pipelines, some open-source, are available to manage RAD data, thereby mitigating  
56 the costs of data analysis. Consequently, RAD sequencing is adept at handling large-scale  
57 projects and proves more cost-effective for the genotyping of numerous samples or a  
58 comprehensive set of SNPs. RAD-seq's appeal extends beyond its economic and temporal  
59 advantages, as it is also adept at analyzing complex genomic sequences, unconstrained by the  
60 size of the genome. This technique has been effectively leveraged in the construction of linkage  
61 maps, comparative genomics, and population genetics studies across a diverse array of  
62 organisms

63

64 Parentage analysis can generally be divided into two categories: exclusion and likelihood  
65 (Flanagan et al. 2019). The former is designed to identify incompatibilities between pairs of  
66 individuals to prove that one cannot be the parent of the other (Chakraborty et al. 1974).  
67 Parentage assignment using likelihood is well-developed approach (Kalinowski et al. 2007;  
68 Marshall et al. 1998). However, different from the rapid development of parentage assignment  
69 methods and related software over the past decade (Flanagan et al. 2019), little work has  
70 focused on generation classification, an important first step.

71

72 Due to its high carcass yield, rapid growth rate, and marbled meat quality, the Angus breed has  
73 become the predominant cattle breed in numerous countries. Xiangxi yellow cattle is an  
74 indigenous Chinese breed and was included in the *National Protection List of Livestock and*

75 *Poultry Genetic Resources of China* in 2006 (Li et al. 2014). This breed can feed on low-quality  
76 roughage, is good at climbing hills, and is well-adapted to high temperatures. To improve meat  
77 quality and quantity, Angus was crossbred with Xiangxi yellow cattle, creating a new hybrid  
78 known as Xiangzhong black cattle, which captures the beneficial traits of each breed. The F1  
79 generation backcrossed with sires to produce the F2 generation. To reduce breeding costs, it is  
80 common to breed F1 and F2 generations in one barn after crossbreeding. However, as the size  
81 of the crossbred population increases, inaccurate recording of the generation and parentage are  
82 more frequent due to error. Therefore, assigning individuals to each generation of the  
83 crossbreeding program is necessary, especially for genetic improvement and local resource  
84 conservation.

85

86 From the F1 generation, the phenotype of the offspring is intermediate between those of parent  
87 and the female, and it is difficult to distinguish the F1 and F2 generations just by inspection. In  
88 the absence of accurate records from local livestock farms, it is nearly impossible to distinguish  
89 between the F1 and F2 generations. Therefore, in this study, we evaluated 975 cattle from  
90 livestock farms by RAD-seq, including sires (Angus), dams (Xiangxi yellow cattle), the F1  
91 generation from a cross between dams and sires, and the F2 generation from the backcross  
92 between F1 and sires and the intercross between F1. Our aims were 1) to distinguish between  
93 the F1 and F2 generations according to Mendelian laws of inheritance and 2) to construct a SNP  
94 panel with high confidence for kinship analysis and pedigree reconstruction for Xiangzhong  
95 black cattle.

## 96 **Material and Methods**

### 97 **Ethics statement**

98 We declare that all of the experimental procedures involving animals were conducted following  
99 the guidelines for experimental animals and were approved by the institutional review board of  
100 BGI (approval number BGI-IRB E22013). The samples were collected from a local farm in  
101 Loudi, Hunan Province, China.

## 102 **Sample information**

103 In total, 975 cattle were sampled, including 923 calves, 28 dams, and 24 sires. First, 28 Xiangxi  
104 yellow cattle dams were crossed with 24 Aberdeen-angus sires. Semen was collected from  
105 Angus bulls and crossed with females of the F1 generation to produce the F2 generation. All  
106 F1 and F2 individuals were cultured together in livestock farms. Mating between bulls and cows  
107 of the F1 generation was possible. Sires and dams were clearly labeled; however, for offspring  
108 cattle, the paternal line was ambiguous.

## 109 **Library preparation and RAD-seq**

110 For DNA extraction, fresh blood was collected and prepared according to the protocol of  
111 solution-based DNA extraction methods (Chacon Cortes et al. 2014). After quality control,  
112 DNA samples were subjected to single-digest RAD-seq as described previously, with slight  
113 modifications (Baird et al. 2008). Briefly, samples meeting quality standards were digested with  
114 TaqI restriction endonuclease for 20 min at 37°C and then randomly fractionated by Covaris  
115 Focused-Ultrasonicator. Then, 5 µl of RAD adapter was added to the interruption product. The  
116 products were cyclized by rolling circle amplification. The libraries were sequenced on the  
117 BGISEQ-500 platform (BGI, Shenzhen, China).

## 118 **SNP genotyping and selection**

119 The raw sequence data were filtered and trimmed using SOAPnuke (Chen et al. 2018) and  
120 separated according to the unambiguous barcodes and the specific enzyme recognition site.  
121 Each RAD read was mapped to the *Bos taurus* genome  
122 ([https://ftp.ncbi.nih.gov/genomes/all/GCF/002/263/795/GCF\\_002263795.1\\_ARS-UCD1.2/](https://ftp.ncbi.nih.gov/genomes/all/GCF/002/263/795/GCF_002263795.1_ARS-UCD1.2/))  
123 using BWA (v0.7.12-r1039) (Li et al. 2010). The bam format was sorted and indexed using  
124 Samtools (v1.14) with default parameters except markdup (Li et al. 2009). Then the format was  
125 transefered to sam format using Picard (v2.26.10). , The GATK (v4.1.2) were used for SNP  
126 and InDel calling with the Variant Call Format (VCF) file with default parameters (McKenna  
127 et al. 2010).

128 Quality control of candidate SNPs was performed using GATK "VariantFiltration." All SNPs  
129 that met the criteria for minimizing false positives ( $QUAL < 250.0$ ,  $DP < 1500$ ,  $DP > 6000$ ,  
130  $MQ < 50.0$  ||  $QD < 4.0$  ||  $FS > 15.0$  ||  $BaseQRankSum < -4.0$  ||  $ReadPosRankSum < -3.50$  ||  
131  $MQRankSum < -10.0$ ,  $SOR > 4.0$ ,  $AN < 975$ ,  $MQ0 \geq 30$ ) were considered potential high-quality  
132 SNP markers for subsequent analyses.

133  
134 For the classification of generations, considering that there were 24 and 28 sires and dams,  
135 respectively, alleles with counts of 48 and 56 indicated homozygous loci. To expand the criteria  
136 to obtain adequate loci, we allowed for five different alleles, which could reflect sequencing  
137 error. Therefore, we selected allele counts of 43–48 and 51–56 in the paternal and maternal  
138 genome along with different allelic genotypes for further analysis. Loci with allele depths below  
139 6 were excluded from the SNP set.

140  
141 For parentage assignment, SNPs with more than two alleles were first excluded. The retained  
142 SNPs from paternal data meeting the following criteria were extracted to form an allele set:  
143 MAF (minor allele frequency)  $> 0.25$ , heterozygosity  $> 0.2$ , proportion of missing data  $< 0.7$ ,

144  $p$ -value for Hardy–Weinberg equilibrium  $> 0.01$ . Maternal and offspring data were filtered  
145 using this allelic set. Finally, sites with an allele depth below 6 and proportion of missing data  
146 above 0.8 were excluded.

## 147 **Classification of F1 and F2 generations**

148 All samples were genotyped based on the SNP panels used for generation classification and  
149 loci that were identical to the maternal/paternal allele or were heterozygous were counted. The  
150 heterozygous ratio was defined as the number of heterozygous SNPs divided by the total SNPs  
151 genotyped. The number of identical maternal/paternal SNPs divided by the total number of  
152 SNPs successfully genotyped was defined as the ratio of SNPs inherited from maternal/paternal  
153 parents. Samples with a heterozygous ratio above 0.5 were classified as F1 and the remaining  
154 samples were assigned to the F2 generation.

155  
156 Scatter diagrams of the results were generated using R. For phylogenic tree construction, the  $p$ -  
157 distance matrix was calculated using VCF2Dis-1.44 ([https://github.com/BGI-](https://github.com/BGI-shenzhen/VCF2Dis)  
158 [shenzhen/VCF2Dis](https://github.com/BGI-shenzhen/VCF2Dis)) and a neighbor-joining tree was generated using PHYLIPNEW with  
159 fneighbor (v3.69.650) (Rice et al. 2000). Finally, visualization was performed using Mega  
160 (Tamura et al. 2007).

161  
162 The population genetic structure was evaluated using the program Admixture. The number of  
163 assumed genetic clusters  $K$  was set to 5, with 10,000 iterations for each run. The SNP selection  
164 was performed based on the following strategies: SNPs were chosen from loci with frequencies  
165 ranging from 43 to 48 and 51 to 56 among paternal and maternal data, which exhibited distinct  
166 allelic genotypes. The SNP panel comprised a total of 7305 loci. Identity-by-descent (IBD) was  
167 analyzed using the PLINK kinship matrix generated based on 7305 SNPs using the kinship

168 matrix function (-method Centered\_IBS) in TASSEL (Bradbury et al. 2007) (v5.0) and  
169 visualized using the heatmap package of R.

170

171

## 172 **Parentage assignment**

173 Parentage assignment was performed using CERVUS (v3.0.7) with the paternity selection  
174 function based on the likelihood approach. Based on a simulation, 959 loci showed 91%  
175 successful assignment at the 95% confidence level. The R package APIS (v1.0.1) (Griot et al.  
176 2020) was used for parentage assignment based on the distribution of Mendelian transmission  
177 probabilities. The assignment error rate was set to 0.05. The first candidate sire given by the  
178 software was considered the most probable sire for the offspring and was used for further  
179 analyses. Exclusion-based parentage assignment was performed using the hiphop package  
180 (v0.0.1) in R (Cockburn et al. 2021). None of the dams and sires were defined as social parents.  
181 The first-ranked candidate parent pair was selected as the most likely parent and subjected to  
182 subsequent analyses.

## 183 **Results**

### 184 **Discovery of SNPs by RAD-seq**

185 To enhance the quality of beef and maintain local genetic resources, local farms introduced  
186 Aberdeen Angus as sires and crossed them with a local breed, Xiangxi yellow cattle (Figure  
187 1A). Female F1 cattle were backcrossed with sires to produce the F2 generation (Figure 1A).  
188 Theoretically, if a sire and dam have different genotypes (AA and aa), the F1 generation will  
189 be heterozygous (Aa) and the F2 generation will contain two genotypes (AA, Aa). The F2



190 generation can be used for further crossbreeding, genetic improvement, and resource  
191 conservation. To reduce costs, farmers usually do not breed F1 and F2 generations in separate  
192 barns. Furthermore, mating is possible between bulls and cows of the F1 generation. Angus  
193 cattle and Xiangzhong Yellow cattle differ mainly in body color and size, i.e., Angus cattle are  
194 black and Xiangzhong Yellow cattle are brown. However, their offspring are mainly black and  
195 it is difficult to distinguish between F2 progeny based on appearance, without genotyping. The  
196 phenotype of F1 hybrids is similar to that of F2 carrying AA and Aa genotypes. In this study,  
197 we evaluated 975 cattle by RAD-seq, including 923 offspring with 24 known sires and 28  
198 known dams. All samples were collected from Hunan province, China.

199

200 A total of 4.65 Tb of clean data were obtained from the samples, with a mean of 35.59 million  
201 reads, 1.59× depth, and 31.06% coverage per individual. An average of 35.96 million reads  
202 were mapped to the *B. taurus* reference genome ARS-UCD1.2, with an average mapping rate  
203 of 99.8%.

204

205 Genetic variants were initially detected by GATK and then underwent preliminary filtering. A  
206 total of 8,155,410 SNPs and 1,100,880 InDel variants that satisfied the criteria were retained  
207 (Supplementary Tables 1,2). A density plot revealed that SNPs and InDels were preferentially  
208 distributed in the proximal telomeric regions (Supplementary Figures 1,2). NC\_037357.1, the  
209 X chromosome of *B. taurus*, had relatively few SNPs and InDels. The X-chromosome is present  
210 in a single copy in males. Compared with those of autosomes, the X-chromosome has a lower  
211 recombination rate, lower mutation rate, and smaller effective population size (Schaffner,  
212 2004).

## 213 **Construction of a SNP panel retaining parental information**

214 Our first aim was to assign cattle lacking clear identifying information to each generation of the  
215 crossbreeding program. We had accurate generation information of 28 offspring from the farm  
216 record. According to Mendel's laws of inheritance, in a dominant-recessive inheritance system,  
217 parents with different homozygous genotypes (AA and aa) will produce heterozygous F1  
218 offspring (Aa). In the backcross of F1 to sires, all F2 progeny show the dominant trait, 50% of  
219 individuals will be homozygous and 50% will be heterozygous, on average. If the F1 are  
220 intercrossed, on average, the F2 generation is expected to be 25% homozygous with the  
221 dominant trait, 50% heterozygous showing the dominant trait (genetic carriers), and 25%  
222 homozygous with the recessive trait. The genotypic ratio is 1 (AA):2 (Aa):1 (aa), and the  
223 phenotypic ratio is 3:1. Accordingly, homozygous loci from parents with different genotypes  
224 may provide pivotal information for identifying the generation of individuals (Figure 1B).

225

226 Retained SNPs were further filtered to ensure the maximum difference between paternal and  
227 maternal genome information. The workflow used to generate a high confidence SNP panel  
228 consisting of a series of the homozygous loci is provided in Figure 1C. First, we chose SNPs  
229 that were homozygous in sires (AA) and dams (aa) for alternative alleles. However, only 45  
230 loci were kept from 24 paternal and 28 maternal SNPs. Considering that the sequencing depth  
231 of some samples was not sufficient for genotyping, the criteria were expanded slightly to obtain  
232 adequate loci. The loci with frequencies of 43–48 and 51–56 among paternal and maternal data  
233 and exhibiting different allelic genotypes were selected as candidates for classification. Finally,  
234 the SNP panel included 7305 loci.

235

236 Subsequently, we constructed a phylogenetic tree based on the panel with data for 24 sires, 28  
237 dams, 15 offspring of the F1 generation, and 13 offspring of the F2 generation. All sires and  
238 dams were clearly assigned to distant groups in the phylogenetic tree. In addition, the F1 and

239 F2 generations were generally separated (Figure 2A). However, it was still difficult to  
240 distinguish between F1 and F2 based on these data. Accordingly, it is necessary to consider  
241 both the percentage of heterozygous loci and the percentage of loci inherited from dams and  
242 sires.

### 243 **Assigning individuals to F1 and F2**

244 From data for the F1 and F2 generation, we extracted the genotypes of all samples using the  
245 high confidence SNP panel described above and compared them with the genotypes of sires  
246 and dams. The numbers of SNP loci in offspring with the maternal/paternal genotype  
247 (paternal/paternal genotype locus, MGL/PGL) were used to calculate the maternal/paternal  
248 ratio, and heterozygous loci (H-genotype) were used to calculate the heterozygous ratio. The  
249 maternal/paternal ratio was defined as the ratio of MGL/PGL to all genotypes. As mentioned  
250 above, all samples belonging to the F1 generation should have heterozygous alleles and these  
251 alleles should segregate independently in the F2 generation. However, since we extended the  
252 criteria for the SNP panel, it is possible to obtain MGL or PGL in the F1 generation.  
253 Consequently, the offspring samples with heterozygous locus percentages above 0.5 were  
254 assigned to the F1 generation and the rest were classified as F2. In the F2 generation, offspring  
255 with a maternal/paternal ratio above 0.5 were defined as having the most maternal/paternal  
256 characteristics.

257

258 As expected, 21 out of 28 offspring were correctly assigned based on farm records. A scatter  
259 diagram of all samples revealed four groups corresponding to sires (uppermost group in Figure  
260 2B), dams (lower), F1 (left corner), and F2 (near dams) (Figure 2B). Five samples were  
261 assigned to the F1 generation but belonged to the F2 generation. It is possible that the F2  
262 individuals were produced from F1 inter-crossing, which influenced the final classification

263 results (Figure 1B). However, considering that the selection of superior individuals with respect  
264 to performance traits from sires and dams is the ultimate goal of this crossbreeding system, it  
265 is not necessary to distinguish between these individuals.

266

267 The newly developed SNP panel was used to genotype a total of 975 samples. Based on SNP  
268 information, almost all offspring were categorized successfully (Supplementary Table 3).  
269 Samples were clearly divided into five groups: the F1 generation (heterozygous ratio  $> 0.5$ , 616  
270 samples), the F2 generation with strong paternal characteristics (MSF2, heterozygous ratio  $>$   
271  $0.5$  and paternal ratio  $> 0.5$ , 227 samples), the F2 generation with strong maternal characteristics  
272 (MDF2, heterozygous ratio  $> 0.5$  and maternal ratio  $> 0.5$ , 73 samples), the F2 generation with  
273 slight paternal characteristics (LSF2, heterozygous ratio  $> 0.5$  and paternal ratio  $>$  maternal  
274 ratio, 37 samples), and the F2 generation with slight maternal characteristics (LDF2,  
275 heterozygous ratio  $> 0.5$ , maternal ratio  $>$  paternal ratio, 26 samples). The scatter diagram  
276 supported the classification of samples into five groups (Figure 2C). Subsequently, we  
277 generated a phylogenetic tree based on the information for all samples (Figure 2D). Samples  
278 classified as F2 were distributed across the phylogenetic tree, with one group clustered near  
279 sires (MSF2) and another group clustered near dams (MDF2). The samples located in an  
280 intermediate position were considered F1 or F2 with slight parental characteristics (LSF2 and  
281 LDF2) (Supplementary Table 3).

## 282 **Validation of classification results**

283 To confirm the correctness of classification, we performed a population structure analysis using  
284 all SNPs to observe the systematic differences in allele frequencies among subpopulations  
285 (Figure 3A). In total, 975 samples were divided into five subpopulations, consistent with the  
286 results for the classification of generations.

287

288 Gene IBD is a fundamental concept that explains genetically mediated similarities among  
289 relatives (Thompson, 2013). The  $PI_{HAT}$  value ( $\hat{\pi}$ , *proportion IBD*) obtained using PLINK  
290 provides an estimate of IBD. This method is based on the hidden Markov model (HMM) and  
291 calculates the probability of  $IBD = 1, 2,$  or  $0$  by method-of-moments estimation. The value of  
292  $PI_{HAT}$  is between  $0$  and  $1$ , and a higher value indicates a closer relationship. We performed  
293 an IBD analysis to evaluate relationships among individuals. Usually, pairs with  $PI_{HAT}$  near  
294  $1$  are considered identical twins, those with values of approximately  $0.5$  are considered  
295 parent/child or parental identical twins, and those with values of around  $0.25$  are considered  
296 half-siblings. We used the homozygous sites from sires and dams to calculate  $PI_{HAT}$ . All  
297 pairs with a  $PI_{HAT}$  value above  $0.5$  were pairs between sires or dams. Most pairs with  
298  $PI_{HAT}$  values around  $0.5$  ( $0.3$ – $0.5$ ) were pairs between sires/dams and MSF2/MDF2  
299 (Supplementary Table 7). These results were in accordance with our previous assignments to  
300 provide a basis for selecting MSF2 for breeding and other F2 for conservation.

301

302 Finally, we created kinship matrices for individuals containing sires, dams, and the F1  
303 generation. These were multiplied by two to indicate expected covariances between samples  
304 and were used for parentage assignment (Figure 3B, Supplementary Table 4). This approach  
305 proved to be effective for distinguishing between heterozygous offspring and candidate  
306 homozygous offspring with parental traits.

### 307 **Parentage assignment of the F1 generation**

308 For parentage assignment with the F1 generation, we chose three tools with different underlying  
309 algorithms and compared the results. CERVUS is widely used and utilizes a likelihood-based

310 approach, APIS uses the observed distribution of Mendelian transmission probabilities, and  
311 HIPHOP extends exclusion approaches with SNP markers.

312

313 In total, 960 SNPs that met the quality criteria were further filtered to determine loci that could  
314 provide sufficient statistical power for parentage assignment. CERVUS requires allele  
315 frequency estimation and a simulation before formal assignment. Simulation can be used to  
316 examine the feasibility of parentage analysis and calculate critical values of likelihood ratios to  
317 determine the confidence of parentage assignments for real data. The following parameters  
318 were used for simulation: 30% loci typed, 0.01% genotyping error rate, 10,000 offspring, 50  
319 candidate fathers, 90% of candidate fathers sampled, and the minimum number of typed loci  
320 set to 10. The output of strict confidence value of simulation was 95%. Finally, by input the  
321 simulation and SNP files, the set containing 28 sires and the panel of 960 SNP markers provided  
322 99.5% successful parentage assignment for 615 offspring (Supplementary Table 5). Among 612  
323 F1–sire relationships, 552 had confidence levels of >80%, including 451 with confidence  
324  $\geq 95\%$ , which indicated a highly significant relationship. The accepted assignment error rate  
325 was 0.5 using APIS (Supplementary Figure 3). HIPHOP required a known social parent and  
326 year of birth for each cattle at the time of parentage assignment. If the individual is the social  
327 parent of the brood, then the social parent parameter was set to 1; otherwise, the parameter was  
328 set to 0. All cattle were bred in livestock farms and none were classified as a social parent. We  
329 then compared the results obtained by the three algorithms and drew a Venn diagram to visually  
330 evaluate shared assignments (Figure 3B). The results obtained using CERVUS and APIS  
331 overlapped by more than 75% (450 identical assignments), while HIPHOP assignments agreed  
332 with those of the other algorithms in <50% of cases (Figure 3C, Supplementary Table 5).

333

334 As mentioned before, HIPHOP relies on exclusion methods. This method identifies  
335 incompatibilities between pairs of individuals according to Mendel's laws, and its accuracy  
336 depends on the accuracy of the marker data. When the power of the marker set is low, the results  
337 become unauthentic. However, it is nearly impossible to obtain error-free data as the sample  
338 size and number of markers grow. Consequently, the overlapping assignments obtained by  
339 CERVUS and APIS were considered the likely parentage assignments with high confidence  
340 (Supplementary Table 6).

## 341 **Discussion**

342 RAD-Seq to genotype SNPs is a practical and cost-effective approach compared to traditional  
343 methods such as microarrays. Choosing informative SNPs reduces the computational burden of  
344 downstream analyses, and allows us to obtain the desired level of accuracy without generating  
345 excessive amounts of data or compromising quality. Overall, RAD-Seq remains an efficient  
346 tool for genotyping and has wide applications in population genetics research. Our SNP panel  
347 offers a cost-effective tool for livestock management by enabling breeders to make informed  
348 decisions on the selection of animals for breeding programs. The detailed genetic information  
349 provided by RAD-Seq can lead to improved health, productivity, and overall welfare of the  
350 livestock. The SNP markers identified in our study can be used to trace and preserve valuable  
351 genetic traits within local breeds. This can facilitate targeted breeding strategies to enhance  
352 specific characteristics, such as disease resistance, adaptability, and production efficiency. This  
353 SNP panel also can be a valuable tool for conservation efforts by providing a means to monitor  
354 and manage the genetic diversity of endangered species. It can also aid in the detection of  
355 hybridization events and inform conservation breeding programs. We propose that future  
356 studies could focus on validating the SNP panel across a broader range of breeds and

357 environments. The SNP panel could be integrated into existing breeding programs, starting with  
358 a pilot study to assess its effectiveness and economic benefits. This would help to establish the  
359 panel's versatility and robustness in different contexts.

360

361 To assign individuals to the F1/F2 generation, a panel of SNPs can be used to determine the  
362 number of loci that are identical to the maternal and paternal alleles or are heterozygous. This  
363 allows for the calculation of maternal, paternal, and heterozygous ratios for each offspring.  
364 Theoretically, the heterozygous ratio for the F1 generation will be around 1.0, while the  
365 maximum ratio was 0.83 in our analysis. Furthermore, the F2 generation was divided into four  
366 subclasses according to the maternal ratio and paternal ratio. This likely reflects the law of  
367 segregation and the law of independent assortment for non-allelic genes on non-homologous  
368 chromosomes, whereas the law of linkage and crossing-over would lead to deviations from the  
369 expected ratio. Unexpectedly, the F2 generation contained the MDF2 and LDF2 subgroups.  
370 Considering that the F1 generation originated from the cross between dams and sires and the  
371 F2 generation from the backcrosses between the F1 generation and sires theoretically, this  
372 observation can likely be explained by crossing within the F1 generation, since the proportions  
373 of MDF2 and LDF2 in the whole F2 generation were low. In addition, during the SNP selection  
374 stage, we allowed five different alleles, which could be the result of sequencing error. This is  
375 another reason why the heterozygous ratio of the F1 generation was less than 1.

376

377 Although several methods for parentage analysis have been developed, accurate parentage  
378 assignment critically depends on the establishment of a molecular marker panel. Accurate  
379 pedigree information has been determined in a Mexican registered Holstein population (García-  
380 Ruiz et al. 2019) and in Chinese Simmental cattle based on a high-density SNP array (Zhang et  
381 al. 2018). Most parentage assignments rely on a microsatellite chip or SNP chip. In taxa without



382 accurate SNP chips for parentage assignment, it is difficult to achieve this goal. RAD-seq and  
383 WGS provide abundant variant information, which enables parentage and population analyses  
384 at the same time. Of note, the SNP panel used for parentage assignment should be optimized  
385 based on several factors, including linkage disequilibrium, minor allele frequencies, deviation  
386 from Hardy-Weinberg equilibrium, genotyping errors, and the frequency of null alleles. In this  
387 study, we achieved parentage assignment by extracting the intersection between CERVUS and  
388 APIS results. However, 165 offspring could not be assigned to accurate sires. This could  
389 probably be attributed to an inherent defect of RAD-seq resulting in a high frequency of null  
390 alleles. Null alleles cannot provide accurate predictions, thus weakening the confidence of the  
391 prediction results.

392

393 Another feasible approach to achieving parentage assignment is based on the genotypes of sex  
394 chromosomes and mitochondrial chromosomes. Typically, the X chromosome and Y  
395 chromosome do not undergo standard recombination in males, and mitochondrial chromosomes  
396 are inherited directionally from the mother (Ballard et al. 2004; Birky Jr 2001; Schaffner 2004).  
397 Therefore, alleles on the sex chromosome and mitochondria will be transferred to offspring  
398 directly. The Y chromosome is difficult to assemble because it contains many ampliconic and  
399 palindromic sequences (Colaco et al. 2018). It is also obvious that SNPs and InDels were  
400 sparsely distributed on the X chromosome (Supplementary Figure 1,2) except for the tip region.  
401 In conclusion, genotyping by RAD-seq is an efficient method for the classification of  
402 generations and parentage assignment. Individuals from different generations identified  
403 through the analysis can be used to accelerate the subsequent breeding process and breeding  
404 conservation.

## 405 **Conclusion**

406 Crossbreeding is a widely used and effective tool to increase genetic diversity within a breed.  
407 However, successful crossbreeding relies on accurate marking and recording of newborn  
408 calves, which is laborious and increases in difficulty as the calf population grows. Traditional  
409 SNP genotyping may be more suitable for smaller projects or when specific SNPs are the focus  
410 of the research. RAD sequencing is generally considered a more cost-effective and scalable  
411 solution for large-scale SNP genotyping projects due to its lower per sample costs and ability  
412 to multiplex many samples simultaneously. The choice between the two methods should be  
413 based on the specific requirements of the project, including the number of samples, the number  
414 of SNPs to be genotyped, and the available resources. In this study, we performed RAD  
415 sequencing and identified the F1 and F2 generation from Angus cattle and Xiangxi yellow cattle  
416 crosses according to Mendelian laws of inheritance and selected a SNP panel with high  
417 confidence for a kinship analysis and pedigree reconstruction. The F1 generation and MSF2  
418 can be applied for breed selection and the LDF2 and MDF2 generation can be maintained for  
419 breed conservation. To the best of our knowledge, this is the first application of a RAD-seq-  
420 based approach for simultaneous generation classification and parentage assignment. The  
421 combination of the efficiency of RNA-seq and advances in kinship analysis is expected to  
422 improve breed management, local resource utilization, and conservation.

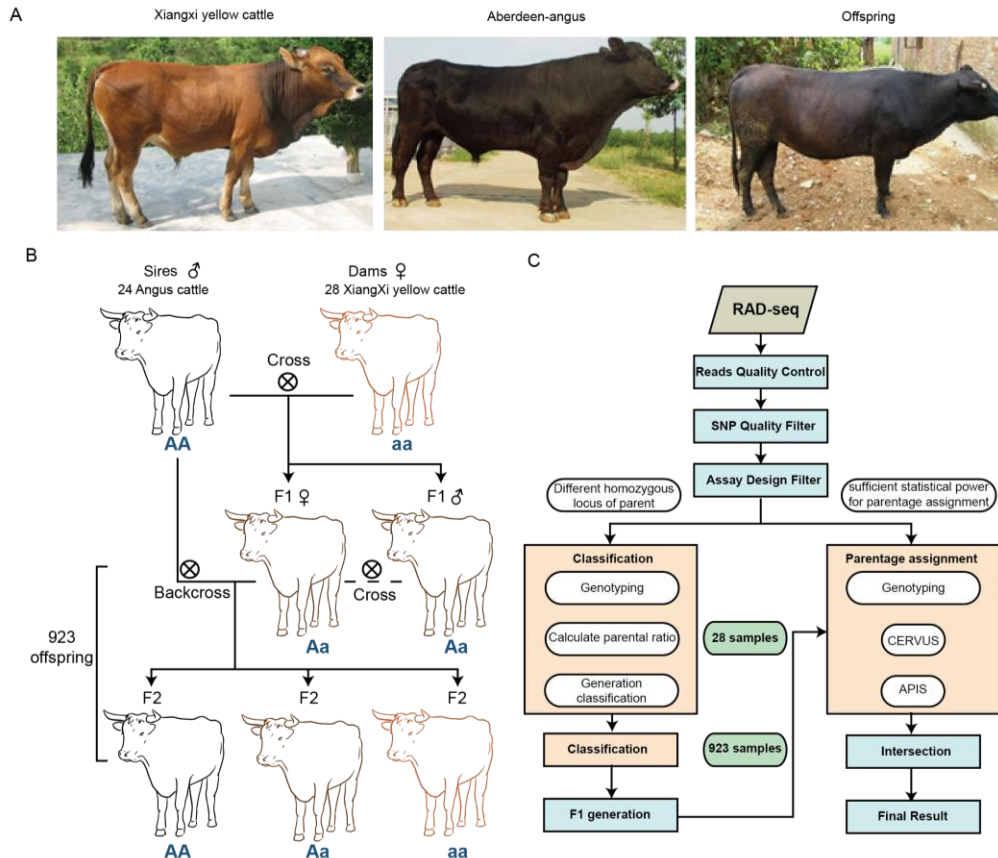
## 423 **Data Availability**

424 The raw sequencing reads were deposited at NCBI under PRJNA1063367, the SNP file was  
425 deposited under figshare (<https://doi.org/10.6084/m9.figshare.25807246>).

426

427

428 **Figures**



429

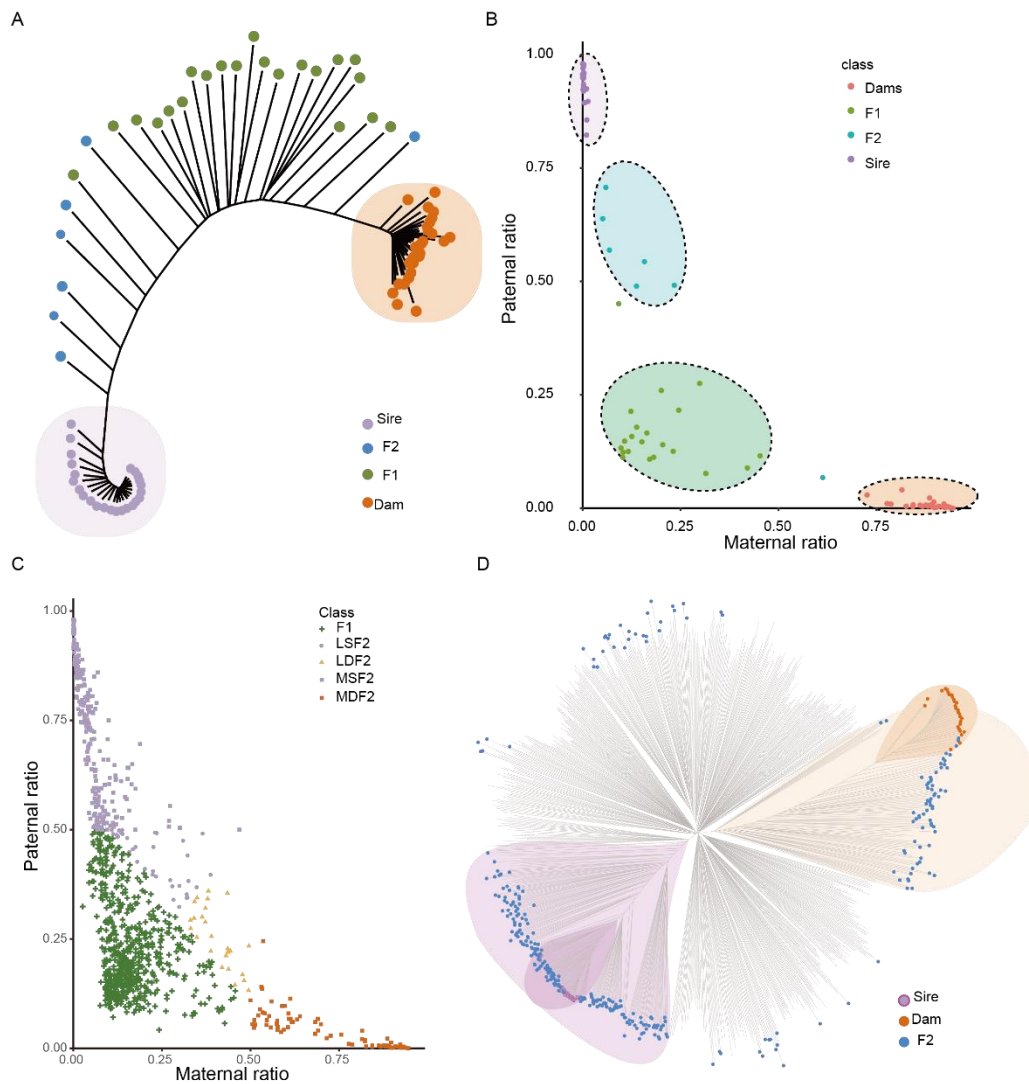
430 **Figure 1: Workflow for generation classification and the selection of a SNP panel**

431 **retaining parental information**

432 A. Representative pictures of the three breeds

433 B. Schematic diagram of the crossbreeding system

434 C. Flowchart detailing the workflow for generation classification and parentage assignment



435

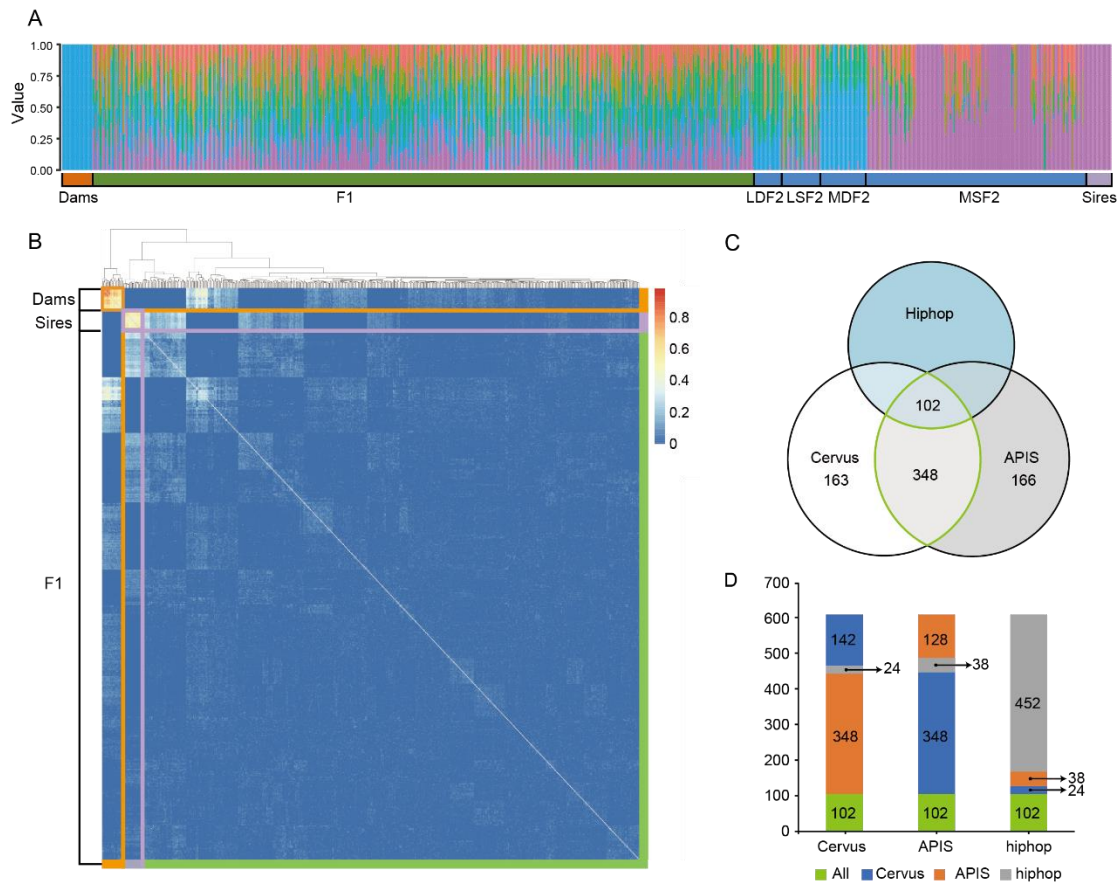
436 **Figure 2: Differentiation between the F1 and F2 generation from a group of cattle**

437 A. Phylogenetic tree of 80 samples containing sires, dams, and offspring with known  
438 generation information.

439 B. Scatter diagram of 80 samples according to the paternal ratio and maternal ratio.

440 C. Scatter diagram of 923 samples according to the paternal ratio and maternal ratio. Different  
441 colors represent different subclasses.

442 D. Phylogenetic tree of all 975 samples containing sires, dams, and offspring.



443

444 **Figure 3: Parentage assignment of the F1 generation**

445 **A.** Population structure analysis of 975 samples with  $K = 5$ .

446 **B.** Heatmap representing the correlations between 615 samples using kinship matrices.

447 **C.** Venn diagram indicating parentage assignments obtained using three algorithms. The  
448 green oval represents the intersection between Cervus and APIS, defined as the final  
449 assignments.

450 **D.** Stacked columns representing the composition within each parentage assignment results  
451 based on three programs.

452

453 **REFERENCES**

- 454 Andolfatto, P., Davison, D., Erezylmaz, D., Hu, T.T., Mast, J., Sunayama-Morita, T., Stern,  
455 D.L., Multiplexed shotgun genotyping for rapid and efficient genetic mapping. *Genome*  
456 *research* 2011. 21, 610–617.
- 457 Baird, N.A., Etter, P.D., Atwood, T.S., Currey, M.C., Shiver, A.L., Lewis, Z.A., Selker, E.U.,  
458 Cresko, W.A., Johnson, E.A., Rapid SNP discovery and genetic mapping using  
459 sequenced RAD markers. *PloS one* 2008. 3, e3376.
- 460 Ballard, J.W.O., Whitlock, M.C., The incomplete natural history of mitochondria. *Molecular*  
461 *ecology* 2004. 13, 729–744.
- 462 Birky Jr, C.W., The inheritance of genes in mitochondria and chloroplasts: laws, mechanisms,  
463 and models. *Annual review of genetics* 2001. 35, 125–148.
- 464 Bradbury, P.J., Zhang, Z., Koon, D.E., Casstevens, T.M., Ramdoss, Y., Buckler, E.S., TASSEL:  
465 software for association mapping of complex traits in diverse samples. *Bioinformatics*  
466 2007. 23, 2633–2635.
- 467 Chacon Cortes, D.F., Griffiths, L., Methods for extracting genomic DNA from whole blood  
468 samples: current perspectives. *Journal of Biorepository Science for Applied Medicine*  
469 2014, 1–9.
- 470 Chakraborty, R., Shaw, M., Schull, W.J., Exclusion of paternity: the current state of the art.  
471 *American Journal of Human Genetics* 1974. 26, 477.
- 472 Chen, Yuxin, Chen, Yongsheng, Shi, C., Huang, Z., Zhang, Y., Li, S., Li, Y., Ye, J., Yu, C., Li,  
473 Z., SOAPnuke: a MapReduce acceleration-supported software for integrated quality  
474 control and preprocessing of high-throughput sequencing data. *Gigascience* 2018. 7,  
475 gix120.
- 476 Cockburn, A., Peñalba, J.V., Jaccoud, D., Kilian, A., Brouwer, L., Double, M.C., Margraf, N.,  
477 Osmond, H.L., Kruuk, L.E.B., Pol, M., HIPHOP : Improved paternity assignment  
478 among close relatives using a simple exclusion method for biallelic markers. *Mol Ecol*  
479 *Resour* 2021. 21, 1850–1865.
- 480 Colaco, S., Modi, D., Genetics of the human Y chromosome and its association with male  
481 infertility. *Reproductive Biology and Endocrinology* 2018. 16, 14.
- 482 Flanagan, S.P., Jones, A.G., The future of parentage analysis: From microsatellites to SNPs and  
483 beyond. *Mol Ecol* 2019. 28, 544–567.
- 484 García-Ruiz, A., Wiggans, G.R., Ruiz-López, F.J., Pedigree verification and parentage  
485 assignment using genomic information in the Mexican Holstein population. *Journal of*



- 486 Dairy Science 2019. 102, 1806–1810.
- 487 Gregory, K.E., Cundiff, L.V., Crossbreeding in Beef Cattle: Evaluation of Systems1. Journal of  
488 Animal Science 1980. 51, 1224–1242.
- 489 Griot, R., Allal, F., Brard-Fudulea, S., Morvezen, R., Haffray, P., Phocas, F., Vandeputte, M.,  
490 APIS: An auto-adaptive parentage inference software that tolerates missing parents.  
491 Mol Ecol Resour 2020. 20, 579–590.
- 492 Kalinowski, S.T., Taper, M.L., Marshall, T.C., Revising how the computer program CERVUS  
493 accommodates genotyping error increases success in paternity assignment. Molecular  
494 ecology 2007. 16, 1099–1106.
- 495 Li, H., Durbin, R., 2010. Fast and accurate long-read alignment with Burrows–Wheeler  
496 transform. Bioinformatics 26, 589–595.
- 497 Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G.,  
498 Durbin, R., The sequence alignment/map format and SAMtools. Bioinformatics 2009.  
499 25, 2078–2079.
- 500 Li, L., Zhu, Y., Wang, X., He, Y., Cao, B., Effects of different dietary energy and protein levels  
501 and sex on growth performance, carcass characteristics and meat quality of F1  
502 Angus × Chinese Xiangxi yellow cattle. Journal of Animal Science and Biotechnology  
503 2014. 5, 21.
- 504 Marshall, T.C., Slate, J., Kruuk, L.E.B., Pemberton, J.M., Statistical confidence for likelihood-  
505 based paternity inference in natural populations. Molecular ecology 1998. 7, 639–655.
- 506 McKenna, A., Hanna, M., Banks, E., Sivachenko, A., Cibulskis, K., Kernytsky, A., Garimella,  
507 K., Altshuler, D., Gabriel, S., Daly, M., The Genome Analysis Toolkit: a MapReduce  
508 framework for analyzing next-generation DNA sequencing data. Genome research 2010.  
509 20, 1297–1303.
- 510 Morin, P.A., Luikart, G., Wayne, R.K., the SNP workshop group, SNPs in ecology, evolution  
511 and conservation. Trends in Ecology & Evolution 2004. 19, 208–216.
- 512 Ng, S.B., Turner, E.H., Robertson, P.D., Flygare, S.D., Bigham, A.W., Lee, C., Shaffer, T., Wong,  
513 M., Bhattacharjee, A., Eichler, E.E., Targeted capture and massively parallel sequencing  
514 of 12 human exomes. Nature 2009. 461, 272–276.
- 515 Rice, P., Longden, I., Bleasby, A., EMBOSS: the European molecular biology open software  
516 suite. Trends in genetics 2000. 16, 276–277.
- 517 Schaffner, S.F., The X chromosome in population genetics. Nature Reviews Genetics 2004. 5,  
518 43–51.

- 519 Tamura, K., Dudley, J., Nei, M., Kumar, S., MEGA4: molecular evolutionary genetics analysis  
520 (MEGA) software version 4.0. *Molecular biology and evolution* 2007. 24, 1596–1599.
- 521 Thompson, E.A., Identity by Descent: Variation in Meiosis, Across Genomes, and in  
522 Populations. *Genetics* 2013. 194, 301–326.
- 523 Zhang, T., Guo, L., Shi, M., Xu, L., Chen, Y., Zhang, L., Gao, H., Li, J., Gao, X., Selection and  
524 effectiveness of informative SNPs for paternity in Chinese Simmental cattle based on a  
525 high-density SNP array. *Gene* 2018. 673, 211–216.
- 526