

## Accelerated Deciphering of the Genetic Architecture of Agricultural Economic Traits in Pigs Using the Low Coverage Whole-genome Sequencing Strategy --Manuscript Draft--

<b>Manuscript Number:</b>	GIGA-D-20-00354	
<b>Full Title:</b>	Accelerated Deciphering of the Genetic Architecture of Agricultural Economic Traits in Pigs Using the Low Coverage Whole-genome Sequencing Strategy	
<b>Article Type:</b>	Research	
<b>Funding Information:</b>	National Transgenic Grand Project (2016ZX08009003-006)	Prof. Xiaoxiang Hu
	948 Program of the Ministry of Agriculture of China (2012-G1(4))	Prof. Xiaoxiang Hu
	Science and Technology Innovation Strategy Projects of Guangdong Province (2019B020203002)	Prof. Xiaoxiang Hu
	Guangdong Academician Workstation (2011A090700016)	Prof. Ning Li
<b>Abstract:</b>	<p>Background: Uncovering the genetic architecture of economic traits in pigs is important for agricultural breeding. Two difficulties limiting the genetic analysis of complex traits are the unavailability of high-density markers for large population in most agricultural species which are lack of good reference panel, and the association signals tend to be spread across most of the genome, i.e., the infinitesimal model of quantitative traits. Findings: Here, we discovered a Tn5-based highly accurate, cost- and time-efficient, low coverage sequencing (LCS) method to obtain whole genome markers and performed whole-genome sequencing on 2,869 Duroc boars at an average depth of 0.73× to identify 11.3 M SNPs. Based on these SNPs, the genome-wide association study (GWAS) detected 14 quantitative trait loci (QTLs) in 7 of 21 important agricultural traits in pigs and provided a starting point for further investigation such as ABCD4 for total teat number and HMGA1 for back fat thickness. The inheritance models of different traits were found to vary greatly. Most obey the minor-polygene model but can be attributed to different reasons, such as the shaping of genetic architecture by artificial selection for this population and sufficiently interconnected minor gene regulatory networks. Conclusions: GWAS results for 21 important agricultural traits identified tens of important QTLs/genes and showed their various genetic architectures, providing promising guidance for genetic improvement harnessing genomic feature. The Tn5-based LCS method can be applied to large-scale genome studies for any species without good reference panel and widely used for agricultural breeding.</p>	
<b>Corresponding Author:</b>	Xiaoxiang Hu China Agricultural University CHINA	
<b>Corresponding Author Secondary Information:</b>		
<b>Corresponding Author's Institution:</b>	China Agricultural University	
<b>Corresponding Author's Secondary Institution:</b>		
<b>First Author:</b>	Ruifei Yang, Ph.D.	
<b>First Author Secondary Information:</b>		
<b>Order of Authors:</b>	Ruifei Yang, Ph.D.	
	Xiaoli Guo	
	Di Zhu	
	Cheng Tan, Ph.D.	

	Cheng Bian
	Jiangli Ren
	Zhuolin Huang
	Yiqiang Zhao, Ph.D.
	Gengyuan Cai, Ph.D.
	Dewu Liu, Ph.D.
	Zhenfang Wu, Ph.D.
	Yuzhe Wang, Ph.D.
	Ning Li, Ph.D.
	Xiaoxiang Hu, Ph.D.
<b>Order of Authors Secondary Information:</b>	
<b>Additional Information:</b>	
<b>Question</b>	<b>Response</b>
Are you submitting this manuscript to a special series or article collection?	No
<p><b>Experimental design and statistics</b></p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	Yes
<p><b>Resources</b></p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite <a href="#">Research Resource Identifiers</a> (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p>	Yes

<p><b>Availability of data and materials</b></p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in <a href="#">publicly available repositories</a> (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our <a href="#">Minimum Standards Reporting Checklist?</a></p>	<p>Yes</p>

1 **Accelerated Deciphering of the Genetic Architecture of**  
2 **Agricultural Economic Traits in Pigs Using the Low Coverage**  
3 **Whole-genome Sequencing Strategy**

4 Ruifei Yang<sup>1,2†</sup>, Xiaoli Guo<sup>1†</sup>, Di Zhu<sup>1†</sup>, Cheng Tan<sup>3</sup>, Cheng Bian<sup>1</sup>, Jiangli Ren<sup>1</sup>, Zhuolin Huang<sup>1</sup>,  
5 Yiqiang Zhao<sup>1</sup>, Gengyuan Cai<sup>3</sup>, Dewu Liu<sup>3</sup>, Zhenfang Wu<sup>3\*</sup>, Yuzhe Wang<sup>1,2\*</sup>, Ning Li<sup>1</sup> and  
6 Xiaoxiang Hu<sup>1\*</sup>

7 <sup>1</sup>State Key Laboratory of Agrobiotechnology, College of Biological Sciences, China  
8 Agricultural University, Beijing, China.

9 <sup>2</sup>College of Animal Science and Technology, China Agricultural University, Beijing, China.

10 <sup>3</sup>National Engineering Research Center for Breeding Swine Industry, South China Agricultural  
11 University, Guangdong, China.

12 Correspondence address: State Key Laboratory of Agrobiotechnology, College of Biological  
13 Sciences, China Agricultural University, Beijing, 100193, P. R. China. Tel: ++86-010-  
14 62733394; E-mail: [huxx@cau.edu.cn](mailto:huxx@cau.edu.cn), [yuzhe891@163.com](mailto:yuzhe891@163.com), [wzfemail@163.com](mailto:wzfemail@163.com).

15 †These authors have contributed equally and should be considered co-first authors

16 **Abstract**

17 Background: Uncovering the genetic architecture of economic traits in pigs is important  
18 for agricultural breeding. Two difficulties limiting the genetic analysis of complex traits  
19 are the unavailability of high-density markers for large population in most agricultural  
20 species which are lack of good reference panel, and the association signals tend to be  
21 spread across most of the genome, i.e., the infinitesimal model of quantitative traits.

22 Findings: Here, we discovered a Tn5-based highly accurate, cost- and time-efficient,  
23 low coverage sequencing (LCS) method to obtain whole genome markers and  
24 performed whole-genome sequencing on 2,869 Duroc boars at an average depth of  
25 0.73× to identify 11.3 M SNPs. Based on these SNPs, the genome-wide association  
26 study (GWAS) detected 14 quantitative trait loci (QTLs) in 7 of 21 important  
27 agricultural traits in pigs and provided a starting point for further investigation such as  
28 *ABCD4* for total teat number and *HMGAI* for back fat thickness. The inheritance  
29 models of different traits were found to vary greatly. Most obey the minor-polygene  
30 model but can be attributed to different reasons, such as the shaping of genetic  
31 architecture by artificial selection for this population and sufficiently interconnected  
32 minor gene regulatory networks. Conclusions: GWAS results for 21 important  
33 agricultural traits identified tens of important QTLs/genes and showed their various  
34 genetic architectures, providing promising guidance for genetic improvement  
35 harnessing genomic feature. The Tn5-based LCS method can be applied to large-scale  
36 genome studies for any species without good reference panel and widely used for  
37 agricultural breeding.

38

39 **KEYWORDS:** Low coverage sequencing; GWAS; Genotyping; Pig; Genetic  
40 architecture; Agricultural traits

41

## 42 **Introduction**

43 Genome-wide association studies (GWAS) have identified thousands of genetic  
44 variants associated with complex traits in humans and agricultural species [1, 2]. The  
45 mapping resolution lies on the density of genetic markers which perceive linkage  
46 disequilibrium (LD) in sufficiently large populations [3, 4]. Despite the declining cost  
47 of sequencing, it is still expensive for agricultural breeding studies to apply whole-  
48 genome sequencing to all individuals in a large cohort (thousands of levels). In many  
49 scenarios, imputation-based strategies, which impute low-density panels to higher  
50 densities, offer an alternative to systematic genotyping or sequencing [5, 6]. To date,  
51 array-based genotype imputation has been widely used in agricultural species [7, 8].  
52 The imputation accuracy of this strategy crucially depends on the reference panel sizes  
53 and genetic distances between the reference and target populations. However, the  
54 unavailability of large reference panels and array designs for target populations in  
55 agricultural species limits the improvement of array-based genotype imputation [9, 10].  
56 Inaccurate imputations influence the results of follow-up population genetic analyses.

57 In terms of recently developed methods, low-coverage sequencing (LCS) of a large  
58 cohort has been proposed to be more informative than sequencing fewer individuals at  
59 a higher coverage rate [11-13]. Sample sizes and haplotype diversity could be more  
60 critical than sequencing depth in determining the genotype accuracy of most  
61 segregating sites and increasing the power of association studies. Overall, LCS has been  
62 proven to have greater power for trait mapping compared to the array-based genotyping  
63 method in human studies [14]. To date, LCS-based genotype imputation has been  
64 employed in many studies using various populations and genotyping algorithms [15-  
65 19]. In particular, the STITCH imputation algorithm overcomes the barrier of the lack  
66 of good reference panels in non-human species and is even applicable in studies with  
67 extremely low sequencing depths [15, 20]. This is a promising approach for agricultural  
68 animals without large reference panels and can be used in the areas of functional genetic  
69 mapping and genomic breeding. However, thus far, no reports on this field have been  
70 done.

71 Several large-scale whole-genome sequencing projects have been completed [21].  
72 These projects were designed to identify the underlying mechanisms that drive  
73 hereditary diseases in human as well as for use in genomic selection in the breeding of  
74 agricultural species [22-24]. The infinitesimal model, which describe the inheritance  
75 patterns of human quantitative traits appears to be successful [25, 26]; however, it is  
76 unclear how many genes play important roles in driving different kinds of complex  
77 traits. In addition, artificial selection provides a driving force to make agricultural  
78 species evolve fast, which further brings about the fixation of selection regions and  
79 differentials in the inheritance model. This process might process a very different result  
80 for the same trait between studies due to different genetic backgrounds of the research  
81 population. Therefore, care should be taken when determining the GWAS result for a  
82 specific population. Such information which might be helpful for understanding the  
83 genetic mechanism for a complex trait and could be informative for further application  
84 of genomic selection in animal breeding.

85 In this study, we developed a new highly accurate, cost- and time-efficient LCS  
86 method to obtain high-density SNP markers for a large Duroc population [27]. By  
87 assessing 21 important agricultural traits in commercial pig herds, we performed  
88 genome-wide association and fine-mapping analyses with high resolution and  
89 compared the results of the inheritance model in depth. We also proved that artificial  
90 selection plays a significant role in altering the genetic architecture of agricultural  
91 animals, especially for those loci that affect economic traits. The LCS strategy provides  
92 a powerful method for further agricultural breeding.

## 93 **Results**

### 94 **Genome sequencing and data acquisition**

95 A Tn5-based protocol was used to prepare sequencing libraries of each pig at a low cost  
96 (reagent cost: \$2.60 /library) as described in the Materials and Methods section. The  
97 libraries were sequenced on the Illumina (PE 150 model, 2 libraries) and the BGI  
98 platform (PE 100 model, 28 libraries) (Supplementary Table S1). The results generated

99 by the BGI platform had lower PCR duplicates (2.23%), higher good index reads  
100 (97.10%), and higher genome coverage (98.55%) than the Illumina dataset (10.82%  
101 PCR duplicates, 93.64% good index reads, and 98.50% genome coverage). Overall, the  
102 total output of the 2,869 boars approached 5.32 TB, and the majority (96.74%) of reads  
103 were successfully mapped to the pig reference genome Sscrofa11.1. Each animal was  
104 sequenced at an average of depth of  $0.73 \pm 0.17 \times$ . Moreover, both high depth  
105 resequencing (average  $15.15 \times$ /sample) and SNP Array (GeneSeek Genomic Profiler  
106 Porcine 80K SNP Array, GGP-80) genotyping were done on the selected Duroc core  
107 boars of this population, and the results were used for downstream accuracy evaluation.

### 108 **Processing pipeline of the low-coverage strategy and accuracy evaluation**

109 Traditional standard methods for SNP calling, such as those implemented in GATK and  
110 Samtools, were mainly used in high-depth resequencing methods. However, due to the  
111 low depth of each base, erroneous SNPs and genotypes could be called using such  
112 methods, especially for the GATK HaplotypeCaller algorithm (single sample local de  
113 novo assembly) [28]. In this study, we applied the BaseVar algorithm [29] to identify  
114 polymorphic sites and infer allele frequencies, and we used STITCH [15] to impute  
115 SNPs. We first used chromosome 18 (SSC18) to test the BaseVar-STITCH and GATK  
116 (UnifiedGenotyper)-Beagle algorithms with genotypes from 1,985 pigs. The 37  
117 verified individuals were genotyped by GATK best practice using HaplotypeCaller for  
118  $15.15 \times$  sequencing data (Fig. 1 and Supplementary Table S2). Correlations ( $R^2$ ) [30]  
119 between genotypes and imputed dosages and the genotypic concordance (GC) were  
120 calculated to evaluate the genotyping accuracy. The initial screening of SSC18 with  
121 BaseVar identified 506,452 and 414,160 bi-allelic candidate polymorphic sites before  
122 and after quality control, respectively. These sites were imputed using STITCH, and  
123 322,386 SNPs were retained with a high average call rate ( $98.89\% \pm 0.59\%$ ) after  
124 quality control (imputation info score  $> 0.4$  and Hardy Weinberg Equilibrium  $P$  value  $>$   
125  $1e^{-6}$ ). The SNPs detected by BaseVar/STITCH were mostly included (99.32%) in the  
126 GATK-Beagle set, which included 570,919 sites and contained 320,199 SNPs  
127 overlapping with the BaseVar/STITCH dataset. As a result, a relatively high-quality



128 genotype set was acquired with less time consumption when  $K = 10$  (the number of  
129 founders or ancestral haplotypes, Fig. S1). Fig. 2 shows that highly accurate genotypes  
130 were obtained using the BaseVar-STITCH pipeline ( $R^2 = 0.919$  and  $GC = 0.970$ ) across  
131 all allele frequencies, which excelled far beyond the method using GATK-Beagle ( $R^2$   
132  $= 0.484$  and  $GC = 0.709$ ). Moreover, we also compared BaseVar-STITCH results with  
133 the genotypes in GGP-80. The results showed even higher GC concordance and  $R^2$   
134 values ( $R^2 = 0.997$  and  $GC = 0.990$ ) when all 2,797 samples were used, which further  
135 validated BaseVar-STITCH with a high level of confidence. Therefore, we conclude  
136 that the BaseVar-STITCH pipeline is a suitable variant discovery and imputation  
137 method for the LCS strategy (Fig. 1).

138 Previous studies have demonstrated that sequencing a large number of samples at a  
139 low depth generally provides a better representation of population genetic variations  
140 compared to sequencing a limited number of individuals at a higher depth. Here, we  
141 examined the consequences of altering the sample size and sequence coverage in this  
142 population. For the  $0.5\times$  coverage using STITCH, a sample size above 500 had little  
143 impact on performance, while at an  $0.1\times$  down-sampled coverage, increasing the  
144 sample size to 1,985 led to a substantially improved performance (Fig. 2C and 2D). At  
145  $0.2\times$  for 1,000 individuals, it was noteworthy that the results were only marginally  
146 poorer ( $R^2 = 0.908$  and  $GC = 0.962$ ) than using all sequencing data (Fig. 2C and 2D).  
147 In general, the total sequencing depth (population category) for one locus  $> 200\times$  was  
148 shown to guarantee the credibility of genotyping within the scope of this study, although  
149 the results did consistently improve as the sequencing depth/sample size increased.

## 150 **Genetic architecture of the Duroc population**

151 After strict parameter filtering in the pipeline (BaseVar-STITCH, Fig. 1), we retained  
152 11,348,460 SNPs in all 2,797 Duroc pigs with high genotype accuracy, and the density  
153 corresponded to 1 SNP per 200 bp in the pig genome (Fig. 3A and Supplementary Table  
154 S3). Finally, the majority of identified SNPs were located in intergenic regions (51.98%)  
155 and intronic regions (36.85%). The exonic regions contained 1.37% of the SNPs,  
156 including 0.14% missense SNPs. Among the discovered SNPs, 1,524,015 (accounting

157 for 13.43% of all SNPs) were novel to the pig dbSNP database (data from NCBI:  
158 GCA\_000003025.6 on June, 2017). Both novel and known variants were found to have  
159 very similar minor allele frequency distributions across the whole genome and the  
160 average minor allele frequency (MAF) was 0.225 (Fig. 3B). A principal component  
161 analysis (PCA) of all pigs showed that there was no distinct population stratification  
162 (Fig. 3C). The decay of LD with increasing distance was different among the  
163 chromosomes, of which the fastest and slowest decay rates occurred for SSC10 and  
164 SSC6, respectively. Average pairwise LD  $r^2$  values fell to 0.20 at 500 Kb and to 0.14 at  
165 1 Mb (Fig. 3D), providing an indication of the expected mapping resolution obtainable  
166 with this population.

167 We further studied the high level of LD and found that it could be a consequence of  
168 long-term strong natural or artificial selection. Tajima's  $D$  and diversity  $\pi$  was  
169 implemented to analyze selective sweep regions simultaneously and only windows with  
170 an interquartile range of Tajima's  $D$  and diversity  $\pi$  of 1.5-fold in the whole genome  
171 were regarded as putative selection regions. In total, 24 putative fixed selective regions  
172 harboring 281 genes were obtained (Fig. S2). The regions displayed significant  
173 overrepresentation of genes involved in the sensory perception of smell ( $P = 6.41e^{-10}$ )  
174 (Supplementary Table S4), reflecting the importance of smell when scavenging for food  
175 during long periods of environmental adaptation. This result is consistent with a  
176 previous study that reported that genes associated with olfaction exhibit fast evolution  
177 in pigs. We also observed a significant enrichment of genes involved in the neurological  
178 system process ( $P = 8.64e^{-5}$ ). These genes may be associated with behavior and  
179 increased tameness and thus were under selection during early domestication. In  
180 addition, the hair cycle process ( $P = 0.004$ ) and bone mineralization ( $P = 0.040$ ) were  
181 also detected to be significantly enriched, which may represent the phenotypic changes  
182 of coat and body composition during pig domestication.

### 183 **GWAS and identification of high-resolution mapping of QTLs**

184 The 21 associated phenotypes used in this study are shown in Table 1 and Fig. S3. We  
185 identified a subset of 258,662 SNPs that tagged all other SNPs with MAF >1% at LD

186  $r^2 < 0.98$  for the first-round of GWAS (Supplementary Table S3). Fine-mapping was  
187 performed within 10 Mb of the SNPs to reach 5% FDR significance threshold genome-  
188 wide. Overall, we discovered a total of 14 non-overlapping QTLs for the seven traits at  
189 a significance threshold of 5% (Fig. 4, Table 1, Fig. S4, and Fig. S5). The widths of all  
190 QTL intervals ranged from ~66 Kb to ~3.9 Mb. The intervals of five QTLs were more  
191 than 2 Mb in width (Supplementary Table S5). These QTLs were strongly influenced  
192 by the local linkage disequilibrium level of this population.

193 On average, each QTL covered 13 protein-coding genes (range of 0–48) with a  
194 median of eight genes. The distribution of the number of genes in a QTL is shown in  
195 Supplementary Table S5. We first focused on QTLs that could be narrowed, since these  
196 loci could provide a starting point for functional investigations. Of the 14 non-  
197 overlapping loci identified in this study, seven QTLs could be further narrowed to a  
198 small number of genes (1 to 9 genes) (Fig. 5 and Fig. S6). Here, we highlight two  
199 important QTLs on SSC7.

200 The QTL on SSC7 with a major effect on the total teat number (TTN) has been  
201 widely identified in several commercial breeding lines and hybrids. Our GWAS results  
202 show a strong QTL for TTN in the same region, explaining most of the phenotypic  
203 variance compared with other QTLs (Supplementary Table S5), reflecting the major  
204 effect of this locus. (Fig. 4). Fine-mapping discovered two narrow LD blocks  
205 (SSC7:97.56–97.65 Mb and 98.06–98.10 Mb) containing four candidate genes (*ABCD4*,  
206 *VRTN*, *PROX2*, and *DLST*) (Fig. 5 and Fig. S6). It is worth noting that four missense  
207 variants were discovered in *PROX2*, one of which was the vertebrate homolog of the  
208 *Drosophila melanogaster* homeodomain-containing protein Prospero, which may be  
209 involved in the determination of cell fate and the establishment of the body plan [31],  
210 and former studies reported that *PROX2* could be the causal gene. Besides, although  
211 there is no direct evidence supporting the involvement of *ABCD4* in the development  
212 process of the mammary gland, we noticed that the most significant locus  
213 (SSC7:97,581,669,  $P = 3.29e^{-22}$ ) was detected in the region of this gene, suggesting that  
214 *ABCD4* may be the most likely causal gene.

215 For the carcass traits, we identified six QTLs (Table 1 and Supplementary Table  
216 S5), in which a common narrowed QTL region on SSC7 of 30.24–30.52 Mb was  
217 identified to be significantly associated with back fat thickness (BF) and loin muscle  
218 depth (LMD) (Fig. 5 and Fig. S6). Among the QTLs associated with BF and LMD, the  
219 narrowed QTL on SSC7 was found to make the greatest contribution to the heritability,  
220 so this would be the location of the major genes in the region (Table 1 and Fig. 5). In  
221 this region (Supplementary Table S5), *HMGAI* is a promising candidate gene  
222 associated with growth, carcass, organ weights, and fat metabolism, as it has been  
223 reported to be involved in a variety of genetic pathways regulating cell growth and  
224 differentiation, glucose uptake, and white and brown adipogenesis [32-36].  
225 *Nudt3* belongs to the Nudix hydrolase family, which is involved in diverse metabolic  
226 processes, including the regulation of important signaling nucleotides and their  
227 metabolites. It is an obesity-linked gene that is associated with insulin signaling and  
228 may be another candidate causal gene of BF and LMD. Other genes, including  
229 *PACSINI* and *SPDEF*, have also been reported to be candidate genes with functions in  
230 glutathione metabolism, adipose and muscle tissue development, and lipid metabolism  
231 for LMD.

### 232 **Heritability and pattern of QTL effects**

233 To assess how much of the heritability can be explained by the detected QTLs, we  
234 estimated the effect size of the overall decreased proportion of heritability by using  
235 significant SNPs distributed in these QTLs as fixed effects. As reported in Table 1, we  
236 detected a larger number of contributions to heritability by major QTLs for the teat  
237 number (3.16~8.86), which indicates that the teat number is mainly controlled by a  
238 small number of loci. We also distributed the effects and significance ( $-\text{Log}_{10} P$  value)  
239 of SNPs for all 21 traits. Again, the result showed that TTN had the most discrete  
240 distribution (Fig. 6).

241 We detected six non-overlapping major QTLs for BF, LMD, and LMP, and the  
242 proportion of explained variation by these QTLs reached 1.19–2.40% which is lower  
243 than that for the teat number. The results reveal that although major QTLs are associated

244 with carcass traits, the effect is relatively limited and there could be a larger number of  
245 minor gene effects.

246 Few QTL were detected for other traits, and most of them could be attributed to the  
247 typically small effect sizes of individual mutations, thousands of which contribute to  
248 the total observed genetic variation for a typical complex trait (such as BH, body length  
249 (BL), and cannon bone (CC)). However, two types of interesting genetic architecture  
250 have caught our attention. In terms of the first one, previous studies reported that growth  
251 traits (such as the average daily gain 30-100 kg (ADG100) and age to 100 kg daily  
252 weight (AGE100)) all have medium or high heritability, and several QTLs have been  
253 detected. However, low heritability traits (ADG100:0.187, AGE100: 0.181) with no  
254 significant QTL were detected in this study. To account for this, we hypothesize that  
255 the major QTL effect may be obscured by rare mutations under strong artificial  
256 selection. We searched the candidate loci of growth traits in the pig QTL database  
257 (<https://www.animalgenome.org/cgi-bin/QTLdb/SS/index>) as well as corresponding  
258 previous reports and identified 51 sites associated with growth traits distributed on 18  
259 chromosomes with a low minor allele frequency ( $MAF < 0.05$ ) in our population.  
260 Moreover, 151 previously-reported candidate sites were not identified as SNPs in this  
261 study (Supplementary Table S6). We checked the sequencing depths of these sites, all  
262 of which exceeded 2,100×, proving that these sites were completely fixed in our  
263 population with the same alleles as in the reference genome. This result reflects the  
264 long-term artificial selection history of this commercial Duroc population for growth  
265 traits and also explains the lost heritability and major QTLs.

266 Second, for the feed intake traits (including average daily feed intake (ADFI),  
267 number of visits to feeder per day (NVD), time spent eating per day (TPD), time spent  
268 eating per visit (TPV), and feed intake per visit (FPV)), the heritability was at a medium  
269 or high level (Fig. 6) but we did not obtain any significant QTL (except one QTL for  
270 TPD). The results showed the distribution of SNP were scattered across the whole  
271 genome (Fig. S7) and the effects was more even (Fig. 6), suggesting that these traits are  
272 controlled by the regulatory effect of multiple minor genes and may have highly

273 complex interactions. In order to clarify the biological functions of these minor  
274 candidate genes, we combined related genes obtained from the top 100 loci from the  
275 GWAS of the six feed intake traits according to the GWAS analysis. The gene-set  
276 enrichment analysis based on the obtained 281 genes showed that neural development  
277 or neural activity related functions, such as astrocyte differentiation ( $P = 8.61e^{-5}$ ),  
278 cognition ( $P = 0.002$ ), learning ( $P = 0.002$ ), and glial cell differentiation ( $P = 0.003$ ),  
279 were significantly enriched (Fig. S7 and Supplementary Table S7). The KEGG pathway  
280 analysis also showed there were significantly enriched nervous system processes (Fig.  
281 S8 and Supplementary Table S8), including the neurotrophin signaling pathway ( $P =$   
282  $0.015$ ) and the GABAergic synapse ( $P = 0.021$ ). This result shows that pig feeding  
283 behavior involves complex traits that are affected by the regulation of the nervous  
284 system, leading to the stimulation of appetite.

## 285 **Discussion**

286 To our knowledge, we have generated the largest WGS genotyping dataset for the  
287 Duroc population so far. It contains 11 million markers from 2,797 pigs. We expanded  
288 the candidate causal mutations for multiple pig traits and demonstrated the efficacy of  
289 genetic fine-mapping utilizing low-coverage sequencing in animal populations without  
290 reference panels. Further, we compared the heritability and inheritance model of each  
291 trait, providing a starting point for functional investigations. Our study indicates that  
292 the LC method could have widespread usage in high resolution genome-wide  
293 association studies for any genetic or breeding population or even for application in  
294 genomic prediction.

295 This study identified an optimal design, taking into account the imputation  
296 algorithm, the number of samples, and the sequencing depth. The BaseVar-STITCH  
297 pipeline allows the GC to be higher than 0.96 when the sample size is 1000 at a  
298 sequencing depth of  $0.2\times$  ( $200\times$  at the population level) without large reference panels.  
299 This GC value is significantly higher than that found in other studies with small sample  
300 sizes with a high sequencing depth or array-based genotype imputation. We also found

301 that the genotype accuracy is more sensitive to the sample size than the sequencing  
302 depth. In other words, the results demonstrated that low-coverage designs are more  
303 powerful than the deep sequencing of fewer individuals for animal sequencing studies,  
304 since a large sample size can cover all local haplotypes of the study population more  
305 effectively. This method has amazing accuracy, even in large-scale human studies with  
306 the most complex population structure [29], which further shows that a sufficient  
307 sample size will ensure that the method has a broad spectrum of applicability in all  
308 agricultural species or any breeding population. Therefore, using the low-coverage  
309 sequencing strategy, we were able to consider both the high-density SNP map and a  
310 large population.

311       Increasing the marker density has been proposed to have the potential to improve  
312 the power of GWAS and the accuracy of genomic selection (GS) for quantitative traits  
313 [37]. First, the whole-genome low-coverage sequencing data gave the best accuracy for  
314 GWAS, since it can catch more recombinations than SNP chips or target sequencing  
315 methods such as genotyping by sequencing (GBS), and most causal or causal-linked  
316 mutations that underlie a trait are expected to be included. Second, lots of studies have  
317 reported the impact of whole genome sequence data on the accuracy of genomic  
318 predictions [37-39]; however, the conclusions have been quite divergent. The limited  
319 improvement of the genetic relationship matrices for WGS data compared with the SNP  
320 chip is the major reason for the lack of improvement in genomic prediction. In addition,  
321 most researchers may prefer to impute SNP chip genotypes using limited WGS data;  
322 however, some erroneous SNPs may be introduced and further adversely affect the  
323 performance of genomic prediction, since limited haplotype architecture would be  
324 obtained using small-scale WGS data. Our method improves the accuracy of imputation,  
325 especially in a large studies without a good reference panel and multibreed genomic  
326 predictions, which will make the application of genome selection wider. Third,  
327 significantly improved GS results have been observed when SNPs were preselected  
328 from the sequenced data using GWAS and a nonlinear genomic prediction method (*e.g.*,  
329 Bayes model [40] or TABLUP [41]). Thus, we could select different useful tag-SNPs

330 for various traits with different genetic architectures using the high-density genetic map  
331 built by LCS data to optimize the genomic selection model in the future. Fourth, in  
332 practical application, the haplotype reference panel can accommodate new haplotypes  
333 due to recombination at any time, thus solving the issue of a decrease in prediction  
334 accuracy over generations. Our data can cover the sites of various SNP chips well  
335 because the genome coverage exceeds 98.36%, and it is competitive with arrays in  
336 terms of the cost and SNP density. Last, we applied GTX, which is an FPGA-based  
337 hardware accelerator platform [42], to do the alignments, and ~3,000 alignments were  
338 accomplished in two days. Then, the genotyping and imputation could be achieved on  
339 the cluster server or even a cloud server in a single day, thus resolving the accuracy and  
340 timeliness issue for genomic prediction.

341 Previous studies demonstrated that pigs have differentiated into a variety of local  
342 populations due to environmental adaptation, and they were domesticated around  
343 ~10,000 years ago. Since then, natural and artificial selection have both contributed to  
344 the further speciation of pigs [43, 44]. Recent swine breeding has prompted the  
345 accumulation of beneficial genetic variations at a more rapid rate, especially for some  
346 economic trait loci [45, 46]. The purebred Duroc population studied in this research  
347 was selected for meat production mainly due to its growth-related index. A large  
348 number of fixed loci have been found to be associated with ADG, AGE, and FCR,  
349 which reflects this selection process exactly. We also detected 24 putative fixed  
350 selective regions. For example, in these regions, *MC5R* was detected to be a possible  
351 candidate gene for fatness in pigs [47], major QTLs for pig growth and carcass traits  
352 were identified to be centered in the regions of *OGN* and *ASPN* [48], and *AKIRIN1* was  
353 found to be involved in the regulation of muscle development by playing important  
354 roles in maintaining the muscle fiber type and regulating skeletal muscle metabolic  
355 activity [49]. *CRTC3* encodes a member of the CRTC protein family and plays an  
356 important role in energy metabolism [50, 51]. It was found to be associated with lipid  
357 accumulation in pigs [51]. In addition, a series of genes enriched for sensory perception  
358 and neurological system processes were also detected in selective regions. It has been



359 widely reported that olfactory receptor genes may not only reflect adaptation to  
360 different environments [43] but also might have acted as a species barrier by affecting  
361 mate choice [52]. Several studies have reported an overrepresentation of genes with GO  
362 (gene ontology) terms related to neuronal development and neurological regulation [43,  
363 53], and this could be related to the complex genetic background of traits such as  
364 behavior and increased tameness. It should be noted that the results may be due to a  
365 mixture of natural and artificial selection causes. The complex genetic background and  
366 single population analyses may limit the precision of exploration of selection signatures  
367 exploration, so analyses of population genetics in multiple breeds in a large population  
368 and multi-omics may be needed.

369 In this study, we detected 14 non-overlapping QTLs in 7 of 21 traits (Table 1 and  
370 Supplementary Table S5). There were big differences of loci and QTL effects among  
371 these traits, which may represent the inheritance models of different traits, including  
372 phenotypes that are mainly affected by several major genes (teat number) or multiple  
373 minor genes (such as carcass traits). Above all, seven non-overlapped QTLs with  
374 narrowed intervals were identified, which emphasizes the potential for identifying new  
375 mutations in QTLs using the low-coverage sequencing method. Some candidate genes  
376 may reside within these regions. For the teat number, we first focused on the QTL  
377 interval on SSC7 which explained most of the phenotypic variance. It should be noted  
378 that six missense SNPs were identified to be extremely significant (Supplementary  
379 Table S9). We estimated the effects of the variants and found one located in *ABCD4*  
380 had the most severe impact with the largest decrease of protein stability. Moreover, the  
381 most extremely significant locus was located in the intron region of *ABCD4* ( $P = 3.29e^{-22}$ ).  
382 We therefore suggest that *ABCD4* is one of the most promising causal genes  
383 affecting the teat number in pigs. For the carcass traits (LMD and BF), several candidate  
384 genes were detected in the narrowed QTLs, especially for the QTL with major effects  
385 on SSC7, including *Grm4*, *Hmgal1*, *NUDT3*, *RPS10*, *PACSINI*, and *SPDEF*, which  
386 have also been widely identified. Moreover, there were three QTLs that had not been  
387 identified in previous studies as far as we know: those detected in TTN

388 (SSC1:34,657,653-36,881,340), LMP (SSC13:83,054,253-84,673,400), and TPD  
389 (SSC1:157,891,084-161,827,351). For TTN and LMP, the newly discovered QTLs  
390 explained the limited phenotypic variance, indicating the minor effects of these loci.  
391 For TPD, we noted that the same QTL was also identified in BF, which suggests that  
392 the intervals may contain genes that control appetite. Apart from the QTLs identified in  
393 high-resolution discussed above, we also detected several loci, though their intervals  
394 could not be narrowed further based on the LD information. Several candidate genes  
395 may affect the regulation or development process which may be worth researching  
396 further (detail in Supplementary Table S10).

397 The GWAS results indicate that the Duroc population delivers fewer loci for fewer  
398 phenotypes. We conclude that the low yield of QTLs can be predominantly explained  
399 by the fixed QTLs for growth traits caused by artificial selection and the infinitesimal  
400 model for high heritability but the lack of major QTL traits. This result shows that the  
401 breeding of this commercial population has been successful, especially in terms of the  
402 improvement of growth traits. The next stage should focus on the use of genomic  
403 selection strategies for “infinitesimal traits” with high heritability but no major QTL,  
404 such as feeding behavior traits. We note that the feeding behavior traits had high or  
405 moderate heritability (Table 1) but a flat SNP distribution compared with TTN (Fig. 6),  
406 which [54] suggests that these traits may rely on a highly polygenic and complex  
407 genetic architecture. According to the GO and KEGG enrichment analyses, mostly  
408 neural activity process related functions or pathways were found to be enriched,  
409 especially the neurotrophin signaling pathway ( $P = 0.015$ ) (Fig. S9). For example, NT3  
410 and TrkB were reported to be involved in the regulation of the nervous system, affecting  
411 the stimulation of appetite [54, 55]. In all, we compared the inheritance models of 21  
412 traits, and the results showed the difference between traits mainly affected by a limited  
413 number of loci and those affected by multiple loci with a small, widely distributed effect.  
414 Moreover, human selection could be a determining factor for the inheritance models of  
415 some specific traits in a specific population, making their genetic mechanism more  
416 complex. For further application of genomic selection, based on the QTL effect and the

417 inheritance model, a suitable prediction model should be designed for breeding to  
418 improve and optimize the accuracy of genomic prediction in animal breeding.

## 419 **Conclusions**

420 In conclusion, we discovered a Tn5-based, highly accurate, cost- and time-efficient  
421 LCS method to obtain whole genome SNP markers in a large Duroc population. We  
422 expect that our method could be applied to large-scale genome studies for any species  
423 without a good reference panel. GWAS results for 21 important agricultural traits  
424 identified tens of important QTLs/genes and showed their various genetic architectures,  
425 providing promising guidance for further genetic improvement harnessing genomic  
426 feature.

## 427 **Methods**

### 428 **Animals, phenotyping, and DNA Extraction**

429 The Duroc boars used for this study were born from September 2011 to September 2013.  
430 All boars were managed on a single nucleus farm in a commercial company, which  
431 enduring strong artificial selection for many years. The associated phenotype data used  
432 in this study included back fat thickness at 100 kg (BF), loin muscle area at 100 kg  
433 (LMA), loin muscle depth at 100 kg (LMD), lean meat percentage at 100 kg (LMP),  
434 average daily gain (0-30 kg and 30-100 kg) (ADG30 and ADG100), age to 30 kg and  
435 100 kg daily weight (AGE30 and AGE100), body length (BL), body height (BH),  
436 circumference of cannon bone (CC), feed conversion ratio (FCR), average daily feed  
437 intake (ADFI), number of visits to feeder per day (NVD), time spent to eat per day  
438 (TPD), time spent to eat per visit (TPV), feed intake per visit (FPV), feed intake rate  
439 (FR), left teat number (LTN), right teat number (RTN), and total teat number (TTN).  
440 The phenotype TTN data were acquired from Tan's study [27]. In detail, the number of  
441 left and right teats of each pig were recorded within 48 h after birth, and only normal  
442 teats were counted. The total teat number in this study was the sum of normal left and  
443 right teats. Body weights were recorded at birth and at the beginning ( $30 \pm 5$  Kg) and

444 the end ( $100 \pm 5$  Kg) of the experiment. The ADG was calculated as the total weight  
445 gain over this time, divided by the number of days. The ages at which the pig reached  
446 30 Kg and 100 Kg were recorded as AGE30 and AGE100 respectively. BF, LMD, LMA,  
447 and LMP were measured over the last three to four ribs using b-ultrasound-scan  
448 equipment when the weight of pigs reached  $100 \pm 5$  Kg (Aloka SSD-500). Feeding  
449 behaviors including the time taken, duration, feed consumption, and weight of each pig  
450 were recorded at every visit by the Osborne FIRE Pig Performance Testing System  
451 (Kansas, American). The ADFI of each animal was obtained by dividing the total feed  
452 intake during the test by the number of days of the test period. The following feeding  
453 behavior and eating efficiency traits were defined and calculated for each boar: ADFI  
454 (Kg/day), TPD (min), NVD, TPV (= TPD/NVD, %), FPV (Kg), FR (= DFI/TPD,  
455 g/min), and FCR (=ADFI/ADG). The phenotypic values nearly all followed a normal  
456 distribution (Fig. S3).

457 Genomic DNA was extracted from the ear tissue using a DNeasy Blood & Tissue  
458 Kit (Qiagen 69506), assessed using a NanoDrop, and checked in 1% agarose gel. All  
459 samples were quantified using a Qubit 2.0 Fluorometer and then diluted to 40 ng/ml in  
460 96-well plates.

#### 461 **Tn5 Library generation and sequencing**

462 Equal amounts of Tn5ME-A/Tn5MErev and Tn5ME-B/Tn5MErev were incubated at  
463 72 °C for 2 minutes and then placed on ice immediately. Tn5 (Karolinska Institute,  
464 Sweden) was loaded with Tn5ME-A+rev and Tn5ME-B+rev in 2× Tn5 dialysis buffer  
465 at 25 °C for 2 h. All linker oligonucleotides were the same as in a previous report [56].

466 Tagmentation were carried out at 55 °C for 10 minutes by mixing 4 µl 5×TAPS-  
467 MgCl<sub>2</sub>, 2 µl of dimethylformamide (DMF) (Sigma Aldrich), 1 µl of the Tn5 pre-diluted  
468 to 16.5 ng/µl, 50 ng of DNA, and nuclease-free water. The total volume of the reaction  
469 was 20 µl. Then, 3.5 µl of 0.2% SDS was added, and Tn5 was inactivated for another  
470 10 min at 55 °C.

471 KAPA HiFi HotStart ReadyMix (Roche) was used for PCR amplification. The  
472 primers were designed for MGI sequencers, with the reverse primers containing 96

473 different index adaptors to distinguish individual libraries. The PCR program was as  
474 follows: 9 min at 72 °C, 30 sec at 98 °C, and then 9 cycles of 30 sec at 98 °C, 30 sec at  
475 63 °C, followed by 3 min at 72 °C. The products were quantified by Qubit Fluorometric  
476 Quantitation (Invitrogen) Then, the groups of 96 indexed samples were pooled with  
477 equal amounts.

478 Size selection was performed using AMPure XP beads (Beckmann), with a left side  
479 size selection ratio of 0.55× and a right side size selection ratio of 0.1×. The final  
480 libraries were sequenced on 2 lanes of MGISEQ-2000 to generate 2×100 bp paired-end  
481 reads or on 1 lane of BGISEQ-500 to generate 2×100 bp paired-end reads.

#### 482 **Genotype data obtained using high depth sequencing and SNP chip**

483 We sequenced 37 out of the total 2,869 pigs using the Hiseq X Ten system at a high  
484 depth of 15.15×. GTX by the Genetalks company, a commercially available FPGA-  
485 based hardware accelerator platform, was used in this study for both mapping clean  
486 reads to the Sscrofa11.1 reference genome ([ftp://ftp.ensembl.org/pub/release-99/fasta/sus\\_scrofa/dna/](ftp://ftp.ensembl.org/pub/release-99/fasta/sus_scrofa/dna/)) and variant calling. The alignment process was accelerated by  
487 FPGA implementation of a parallel seed-and-extend approach based on the Smith–  
488 Waterman algorithm, while the variant calling process was accelerated by FPGA  
489 implementation of GATK HaplotypeCaller (PairHMM) [57]. GATK multi-sample best  
490 practice was used to call and genotype SNPs for the 37 pigs, and the SNPs were hard  
491 filtered with a relatively strict option “QD < 10.0 || ReadPosRankSum < -8.0 || FS >  
492 10.0 || MQ<40.0”. The average running time from a fastq file to a bam file was about 3  
493 min for each sample in this study.

495 We also selected 42 individuals who were included in the LCS dataset and  
496 genotyped using the GeneSeek Genomic Profiler Porcine 80K SNP Array and obtained  
497 68,528 SNPs across the whole genome. The genotypes of the sex chromosomes were  
498 excluded from this study, and after quality control (genotype call rate > 0.95), 47,946  
499 SNPs remained. We retained 45,308 SNPs that overlapped with the LCS dataset to  
500 evaluate the genotypes from the LCS strategy.

501 **Low coverage sequencing data analyses**

502 Sequencing reads from the low coverage samples were mapped to the Sscrofa11.1  
503 reference genome using GTX-align, which includes a step that involves marking PCR  
504 duplicates. The indel realignment and base quality recalibration modules in GATK  
505 were applied to realign the reads around indel candidate loci and to recalibrate the base  
506 quality. Variant calling was done using the BaseVar and hard filtered with EAF  $\geq 0.01$   
507 and a depth greater than or equal to 1.5 times the interquartile range. The detailed  
508 BaseVar algorithm that was used to call SNP variants and estimate allele frequency was  
509 described in a previous report [29]. We used STITCH [15] to impute genotype  
510 probabilities for all individuals. The key parameter K (number of ancestral haplotypes)  
511 was decided based on the tests in SSC18. Results were filtered with an imputation info  
512 score  $> 0.4$  and a Hardy Weinberg Equilibrium (HWE)  $P$  value  $> 1e^{-6}$ . After quality  
513 control, 2,797 individuals with genotype data were obtained. Two validation actions  
514 were taken to calculate the accuracy of imputation. The first parameter was genotypic  
515 concordance (GC), which was calculated as the number of correctly-imputed genotypes  
516 divided by the total number of sites. Another parameter was the allele dosage  $R^2$ , which  
517 was described in a previous report [30]. The SNPEff program [58] was used to annotate  
518 the variants.

519 **Population genetics analysis**

520 A subset of 258,662 SNPs that tagged all other SNPs with MAF  $> 1\%$  at LD  $r^2 < 0.98$   
521 and a call rate of  $>95\%$  were retained for downstream analysis. PCA clustering analyses  
522 were performed with GCTA software [59]. The average heterozygosity rate and MAF  
523 were obtained using the vcftools program [60]. Tajima's D [61] and diversity  $P_i$  was  
524 implemented to analyze selective sweep regions simultaneously with the window size  
525 set to 1 Mb, and only windows with an interquartile range for Tajima's D and diversity  
526  $P_i$  of 1.5 fold in the whole genome were regarded as putative selection regions. Gene  
527 Ontology (GO) and Kyoto Encyclopedia of Genes and Genomes (KEGG) pathway  
528 enrichment analyses were performed using the OmicShare tools  
529 (<http://www.omicshare.com/tools>).

530 **Genome-wide association and Heritability estimation**

531 A mixed linear model (MLM) approach was used for the genome-wide association  
532 analyses, as implemented in the GCTA package [59]. The statistical model included the  
533 year and month as discrete covariates. For BF, LMA, LMD, and LMP, the year and  
534 season were included as discrete covariates, and the weights at the beginning and end  
535 of the test were used as quantitative covariates. To correct for multiple testing across  
536 the genome, the FDR correction obtained using FDRtool R package [62] was applied  
537 to determine the genome-wide significance threshold ( $FDR < 0.05$ ). The SNP effect  
538 was estimated using the GREML\_CE program in the GVCBLUP package [63], where  
539 the result was absolved and normalized.

540 Heritability was estimated using a mixed model as follows:

541 
$$\mathbf{y} = \mathbf{X}_b\mathbf{b} + \mathbf{Z}\mathbf{a} + \mathbf{e}$$

542 with  $\text{Var}(\mathbf{y}) = \mathbf{Z}\mathbf{A}_a\mathbf{Z}'\sigma_a^2 + \mathbf{I}\sigma_e^2$ , where  $\mathbf{Z}$  is an incidence matrix allocating phenotypic  
543 observations to each animal;  $\mathbf{b}$  is the vector of the fixed year-month effects for BF,  
544 LMA, LMD, and LMP that also includes the weights at the beginning and end of the  
545 test as covariance;  $\mathbf{X}_b$  is the incidence matrix for  $\mathbf{b}$ ;  $\mathbf{a}$  is the vector of additive values  
546 based on the genotype data;  $\mathbf{A}_a$  is a genomic additive relationship matrix;  $\sigma_a^2$  is the  
547 additive variance; and  $\sigma_e^2$  is the residual variance. Variance components were estimated  
548 by genomic restricted maximum likelihood estimation (GREML) using the  
549 GREML\_CE program in the GVCBLUP package. The additive heritability was defined  
550 as:  $h_a^2 = \sigma_a^2 / (\sigma_a^2 + \sigma_e^2)$ . SNP effects were defined by the GREML\_CE program and  
551 then normalized using R script.

552 The heritability of the detected QTL was estimated as follows:

553 
$$\mathbf{y} = \mathbf{X}'_b\mathbf{b}' + \mathbf{Z}\mathbf{a} + \mathbf{e}$$

554 with  $\text{Var}(\mathbf{y}) = \mathbf{Z}\mathbf{A}_a\mathbf{Z}'\sigma_a^2 + \mathbf{I}\sigma_e^2$ , where  $\mathbf{Z}$  is an incidence matrix allocating phenotypic  
555 observations to each animal;  $\mathbf{b}'$  is the vector of the fixed year-month effects and  
556 significant SNPs identified in the QTL region using GWAS analysis for BF, LMA,  
557 LMD and LMP;  $\mathbf{b}$  also includes the weights at the beginning and end of the test as  
558 covariance;  $\mathbf{X}'_b$  is the incidence matrix for  $\mathbf{b}$ ;  $\mathbf{a}$  is the vector of additive values based on

559 the genotype data;  $\mathbf{A}_a$  is a genomic additive relationship matrix;  $\sigma_a^2$  is the additive  
560 variance; and  $\sigma_e^2$  is the residual variance. The QTL heritability was defined as  $h_{\text{qtl}}^2 =$   
561  $h_a^2 - \sigma_a^2 / (\sigma_a^2 + \sigma_e^2)$ .

## 562 **Functional Consequence of the Missense Mutations associated with TN**

563 The effect of the missense SNPs associated with TN on the stability of pig ABCD4,  
564 PROX2, and DLST proteins was assessed using I-Mutant adaptation 2.0 [64]. A  
565 potential surge or reduction in the DDG was predicted, along with a reliability index  
566 (RI), where the lowest and highest reliability levels were 0 and 10, respectively.

## 567 **Data availability**

568 All of the sequencing raw data in this study have been deposited into NCBI with  
569 accession number PRJNA681437, and the variance data as VCF file will be available  
570 via GIGADB.

## 571 **Abbreviations:**

572 LCS: Low coverage sequencing method; GC: genotypic concordance; TTN: Total teat  
573 number; LTN: Left teat number; RTN: Right teat number; BF: Back fat thickness at  
574 100 Kg; LMD: Loin muscle depth at 100 Kg; LMA: Loin muscle area at 100 Kg; LMP:  
575 Lean meat percentage at 100 Kg; TPD: Time spent to eat per day; ADFI: Average daily  
576 feed intake; NVD: Number of visits to feeder per day; TPV: Time spent to eat per visit;  
577 FR: Feed intake rate; FPV: Feed intake per visit; FCR: Feed conversion rate; ADG30:  
578 Average daily gain (0-30 Kg); AGE30: Age to 30 kg live weight; ADG100: Average  
579 daily gain (30-100 Kg); AGE100: Age to 100 kg live weight; BL: Body length; BH:  
580 Body height; CC: Circumference of cannon bone.

## 581 **Competing interests**

582 The authors declare that they have no competing interests.

## 583 **Funding**



584 This study is supported by the financial support of the National Transgenic Grand  
585 Project [2016ZX08009003-006], the 948 Program of the Ministry of Agriculture of  
586 China (2012-G1(4)), the Science and Technology Innovation Strategy Projects of  
587 Guangdong Province [2019B020203002], and the Guangdong Academician  
588 Workstation [2011A090700016].

### 589 **Authors' contributions**

590 XH, YW, and ZW designed the research. YW and RY led the writing of the paper. RY,  
591 DZ and YW analyzed the data and generated the models. XG, JR, ZH, CB and CT  
592 contributed sequencing method, reagents, materials, and phenotype. RY, YW, XH, YZ,  
593 GC and DL contributed to the interpretation of the results and edited the paper.

### 594 **Acknowledgements**

595 The authors thank Zhuo Song, Chungen Yi, and Wenjuan Wei for providing the  
596 services of FPGA-based hardware accelerator platform and the Batch Compute system  
597 in Aliyun cloud. The authors also thank Siyang Liu and Xun Xu for their valuable  
598 suggestions on data analyses, and Li Ma and Zhaoliang Liu for improving the  
599 manuscript. Part of the analysis was performed on the high-performance computing  
600 platform of the State Key Laboratory of Agrobiotechnology.

### 601 **References**

- 602 1. Visscher PM, Brown MA, McCarthy MI and Yang J. Five years of GWAS  
603 discovery. *Am J Hum Genet.* 2012;90 1:7-24. doi:10.1016/j.ajhg.2011.11.029.
- 604 2. Huang X, Wei X, Sang T, Zhao Q, Feng Q, Zhao Y, et al. Genome-wide  
605 association studies of 14 agronomic traits in rice landraces. *Nat Genet.* 2010;42  
606 11:961-7. doi:10.1038/ng.695.
- 607 3. Marchini J and Howie B. Genotype imputation for genome-wide association  
608 studies. *Nat Rev Genet.* 2010;11 7:499-511. doi:10.1038/nrg2796.
- 609 4. Marchini J, Howie B, Myers S, McVean G and Donnelly P. A new multipoint

- 610 method for genome-wide association studies by imputation of genotypes. *Nat*  
611 *Genet.* 2007;39 7:906-13. doi:10.1038/ng2088.
- 612 5. Howie BN, Donnelly P and Marchini J. A flexible and accurate genotype  
613 imputation method for the next generation of genome-wide association studies.  
614 *PLoS Genet.* 2009;5 6:e1000529. doi:10.1371/journal.pgen.1000529.
- 615 6. Howie B, Fuchsberger C, Stephens M, Marchini J and Abecasis GR. Fast and  
616 accurate genotype imputation in genome-wide association studies through pre-  
617 phasing. *Nature Genetics.* 2012;44 8:955-+. doi:10.1038/ng.2354.
- 618 7. Yan G, Qiao R, Zhang F, Xin W, Xiao S, Huang T, et al. Imputation-Based  
619 Whole-Genome Sequence Association Study Rediscovered the Missing QTL  
620 for Lumbar Number in Sutan Pigs. *Sci Rep.* 2017;7 1:615. doi:10.1038/s41598-  
621 017-00729-0.
- 622 8. van Binsbergen R, Bink MC, Calus MP, van Eeuwijk FA, Hayes BJ, Hulsegge  
623 I, et al. Accuracy of imputation to whole-genome sequence data in Holstein  
624 Friesian cattle. *Genet Sel Evol.* 2014;46:41. doi:10.1186/1297-9686-46-41.
- 625 9. van den Berg S, Vandenplas J, van Eeuwijk FA, Bouwman AC, Lopes MS and  
626 Veerkamp RF. Imputation to whole-genome sequence using multiple pig  
627 populations and its use in genome-wide association studies. *Genetics Selection*  
628 *Evolution.* 2019;51 doi:10.1186/s12711-019-0445-y.
- 629 10. Swarts K, Li HH, Navarro JAR, An D, Romay MC, Hearne S, et al. Novel  
630 Methods to Optimize Genotypic Imputation for Low-Coverage, Next-  
631 Generation Sequence Data in Crop Plants. *Plant Genome-U.S.* 2014;7 3  
632 doi:10.3835/plantgenome2014.05.0023.
- 633 11. Buerkle CA and Gompert Z. Population genomics based on low coverage  
634 sequencing: how low should we go? *Mol Ecol.* 2013;22 11:3028-35.  
635 doi:10.1111/mec.12105.
- 636 12. Huang L, Wang B, Chen RT, Bercovici S and Batzoglou S. Reveel: large-scale  
637 population genotyping using low-coverage sequencing data. *Bioinformatics.*  
638 2016;32 11:1686-96. doi:10.1093/bioinformatics/btv530.

- 639 13. Li Y, Sidore C, Kang HM, Boehnke M and Abecasis GR. Low-coverage  
640 sequencing: implications for design of complex trait association studies.  
641 Genome Res. 2011;21 6:940-51. doi:10.1101/gr.117259.110.
- 642 14. Gilly A, Southam L, Suveges D, Kuchenbaecker K, Moore R, Melloni GEM, et  
643 al. Very low-depth whole-genome sequencing in complex trait association  
644 studies. Bioinformatics. 2019;35 15:2555-61.  
645 doi:10.1093/bioinformatics/bty1032.
- 646 15. Davies RW, Flint J, Myers S and Mott R. Rapid genotype imputation from  
647 sequence without reference panels. Nature Genetics. 2016;48 8:965-+.  
648 doi:10.1038/ng.3594.
- 649 16. Ros-Freixedes R, Gonen S, Gorjanc G and Hickey JM. A method for allocating  
650 low-coverage sequencing resources by targeting haplotypes rather than  
651 individuals. Genetics Selection Evolution. 2017;49 doi:ARTN 78  
652 10.1186/s12711-017-0353-y.
- 653 17. Frago CA, Heffelfinger C, Zhao HY and Dellaporta SL. Imputing Genotypes  
654 in Biallelic Populations from Low-Coverage Sequence Data. Genetics.  
655 2016;202 2:487-+. doi:10.1534/genetics.115.182071.
- 656 18. Bickhart DM, Hutchison JL, Null DJ, VanRaden PM and Cole JB. Reducing  
657 animal sequencing redundancy by preferentially selecting animals with low-  
658 frequency haplotypes. J Dairy Sci. 2016;99 7:5526-34. doi:10.3168/jds.2015-  
659 10347.
- 660 19. Zan Y, Payen T, Lillie M, Honaker CF, Siegel PB and Carlborg O. Genotyping  
661 by low-coverage whole-genome sequencing in intercross pedigrees from  
662 outbred founders: a cost-efficient approach. Genet Sel Evol. 2019;51 1:44.  
663 doi:10.1186/s12711-019-0487-1.
- 664 20. Nicod J, Davies RW, Cai N, Hassett C, Goodstadt L, Cosgrove C, et al. Genome-  
665 wide association of multiple complex traits in outbred mice by ultra-low-  
666 coverage sequencing. Nat Genet. 2016;48 8:912-8. doi:10.1038/ng.3595.
- 667 21. Genomes Project C, Abecasis GR, Altshuler D, Auton A, Brooks LD, Durbin

- 668 RM, et al. A map of human genome variation from population-scale sequencing.  
669 Nature. 2010;467 7319:1061-73. doi:10.1038/nature09534.
- 670 22. Lam HM, Xu X, Liu X, Chen W, Yang G, Wong FL, et al. Resequencing of 31  
671 wild and cultivated soybean genomes identifies patterns of genetic diversity and  
672 selection. Nat Genet. 2010;42 12:1053-9. doi:10.1038/ng.715.
- 673 23. Daetwyler HD, Capitan A, Pausch H, Stothard P, van Binsbergen R, Brondum  
674 RF, et al. Whole-genome sequencing of 234 bulls facilitates mapping of  
675 monogenic and complex traits in cattle. Nat Genet. 2014;46 8:858-65.  
676 doi:10.1038/ng.3034.
- 677 24. Hayes BJ and Daetwyler HD. 1000 Bull Genomes Project to Map Simple and  
678 Complex Genetic Traits in Cattle: Applications and Outcomes. Annu Rev Anim  
679 Biosci. 2019;7:89-102. doi:10.1146/annurev-animal-020518-115024.
- 680 25. Yang J, Benyamin B, McEvoy BP, Gordon S, Henders AK, Nyholt DR, et al.  
681 Common SNPs explain a large proportion of the heritability for human height.  
682 Nat Genet. 2010;42 7:565-9. doi:10.1038/ng.608.
- 683 26. Lango Allen H, Estrada K, Lettre G, Berndt SI, Weedon MN, Rivadeneira F, et  
684 al. Hundreds of variants clustered in genomic loci and biological pathways  
685 affect human height. Nature. 2010;467 7317:832-8. doi:10.1038/nature09410.
- 686 27. Tan C, Wu ZF, Ren JL, Huang ZL, Liu DW, He XY, et al. Genome-wide  
687 association study and accuracy of genomic prediction for teat number in Duroc  
688 pigs using genotyping-by-sequencing. Genetics Selection Evolution. 2017;49  
689 doi:ARTN 35 10.1186/s12711-017-0311-8.
- 690 28. Ros-Freixedes R, Battagin M, Johnsson M, Gorjanc G, Mileham AJ, Rounsley  
691 SD, et al. Impact of index hopping and bias towards the reference allele on  
692 accuracy of genotype calls from low-coverage sequencing. Genet Sel Evol.  
693 2018;50 1:64. doi:10.1186/s12711-018-0436-4.
- 694 29. Liu S, Huang S, Chen F, Zhao L, Yuan Y, Francis SS, et al. Genomic Analyses  
695 from Non-invasive Prenatal Testing Reveal Genetic Associations, Patterns of  
696 Viral Infections, and Chinese Population History. Cell. 2018;175 2:347-59 e14.

- 697 doi:10.1016/j.cell.2018.08.016.
- 698 30. Browning BL and Browning SR. A unified approach to genotype imputation  
699 and haplotype-phase inference for large data sets of trios and unrelated  
700 individuals. *Am J Hum Genet.* 2009;84 2:210-23.  
701 doi:10.1016/j.ajhg.2009.01.005.
- 702 31. Pistocchi A, Bartesaghi S, Cotelli F and Del Giacco L. Identification and  
703 expression pattern of zebrafish prox2 during embryonic development. *Dev Dyn.*  
704 2008;237 12:3916-20. doi:10.1002/dvdy.21798.
- 705 32. Gong H, Xiao S, Li W, Huang T, Huang X, Yan G, et al. Unravelling the genetic  
706 loci for growth and carcass traits in Chinese Bamaxiang pigs based on a 1.4  
707 million SNP array. *J Anim Breed Genet.* 2019;136 1:3-14.  
708 doi:10.1111/jbg.12365.
- 709 33. Liu X, Wang LG, Liang J, Yan H, Zhao KB, Li N, et al. Genome-Wide  
710 Association Study for Certain Carcass Traits and Organ Weights in a Large  
711 WhitexMinzhu Intercross Porcine Population. *J Integr Agr.* 2014;13 12:2721-  
712 30. doi:10.1016/S2095-3119(14)60787-5.
- 713 34. Arce-Cerezo A, Garcia M, Rodriguez-Nuevo A, Crosa-Bonell M, Enguix N,  
714 Pero A, et al. HMGA1 overexpression in adipose tissue impairs adipogenesis  
715 and prevents diet-induced obesity and insulin resistance. *Sci Rep.* 2015;5:14487.  
716 doi:10.1038/srep14487.
- 717 35. Wang LG, Zhang LC, Yan H, Liu X, Li N, Liang J, et al. Genome-Wide  
718 Association Studies Identify the Loci for 5 Exterior Traits in a Large White x  
719 Minzhu Pig Population. *Plos One.* 2014;9 8 doi:ARTN e103766  
720 10.1371/journal.pone.0103766.
- 721 36. Ji JX, Yan GR, Chen D, Xiao SJ, Gao J and Zhang ZY. An association study  
722 using imputed whole-genome sequence data identifies novel significant loci for  
723 growth-related traits in a Duroc x Erhualian F-2 population. *Journal of Animal  
724 Breeding and Genetics.* 2019;136 3:217-28. doi:10.1111/jbg.12389.
- 725 37. Meuwissen T and Goddard M. Accurate prediction of genetic values for

- 726 complex traits by whole-genome resequencing. *Genetics*. 2010;185 2:623-31.  
727 doi:10.1534/genetics.110.116590.
- 728 38. Zhang C, Kemp RA, Stothard P, Wang Z, Boddicker N, Krivushin K, et al.  
729 Genomic evaluation of feed efficiency component traits in Duroc pigs using  
730 80K, 650K and whole-genome sequence variants. *Genet Sel Evol*. 2018;50 1:14.  
731 doi:10.1186/s12711-018-0387-9.
- 732 39. Yan G, Guo T, Xiao S, Zhang F, Xin W, Huang T, et al. Imputation-Based  
733 Whole-Genome Sequence Association Study Reveals Constant and Novel Loci  
734 for Hematological Traits in a Large-Scale Swine F2 Resource Population. *Front*  
735 *Genet*. 2018;9:401. doi:10.3389/fgene.2018.00401.
- 736 40. Meuwissen TH, Hayes BJ and Goddard ME. Prediction of total genetic value  
737 using genome-wide dense marker maps. *Genetics*. 2001;157 4:1819-29.
- 738 41. Zhang Z, Ding X, Liu J, Koning DJD and Zhang Q. Genomic selection for QTL-  
739 MAS data using a trait-specific relationship matrix. *Bmc Proceedings*. 2011;5  
740 Suppl 3 Suppl 3:S15.
- 741 42. Xing Y, Li G, Wang Z, Feng B, Song Z and Wu C. GTZ: a fast compression and  
742 cloud transmission tool optimized for FASTQ files. *BMC Bioinformatics*.  
743 2017;18 Suppl 16:549. doi:10.1186/s12859-017-1973-5.
- 744 43. Groenen MA, Archibald AL, Uenishi H, Tuggle CK, Takeuchi Y, Rothschild  
745 MF, et al. Analyses of pig genomes provide insight into porcine demography  
746 and evolution. *Nature*. 2012;491 7424:393-8. doi:10.1038/nature11622.
- 747 44. Rubin CJ, Megens HJ, Martinez Barrio A, Maqbool K, Sayyab S, Schwochow  
748 D, et al. Strong signatures of selection in the domestic pig genome. *Proc Natl*  
749 *Acad Sci U S A*. 2012;109 48:19529-36. doi:10.1073/pnas.1217149109.
- 750 45. Bosse M, Megens HJ, Frantz LA, Madsen O, Larson G, Paudel Y, et al. Genomic  
751 analysis reveals selection for Asian genes in European pigs following human-  
752 mediated introgression. *Nat Commun*. 2014;5:4392. doi:10.1038/ncomms5392.
- 753 46. Bosse M, Lopes MS, Madsen O, Megens HJ, Crooijmans RP, Frantz LA, et al.  
754 Artificial selection on introduced Asian haplotypes shaped the genetic

- 755 architecture in European commercial pigs. *Proc Biol Sci.* 2015;282  
756 1821:20152019. doi:10.1098/rspb.2015.2019.
- 757 47. Ma Y, Zhang S, Zhang K, Fang C, Xie S, Du X, et al. Genomic Analysis To  
758 Identify Signatures of Artificial Selection and Loci Associated with Important  
759 Economic Traits in Duroc Pigs. *G3 (Bethesda)*. 2018;8 11:3617-25.  
760 doi:10.1534/g3.118.200665.
- 761 48. Stratil A, Van Poucke M, Bartenschlager H, Knoll A, Yerle M, Peelman LJ, et  
762 al. Porcine OGN and ASPN: mapping, polymorphisms and use for quantitative  
763 trait loci identification for growth and carcass traits in a Meishan x Pietrain  
764 intercross. *Anim Genet.* 2006;37 4:415-8. doi:10.1111/j.1365-  
765 2052.2006.01480.x.
- 766 49. Sun W, Hu S, Hu J, Qiu J, Yang S, Hu B, et al. Akirin1 promotes myoblast  
767 differentiation by modulating multiple myoblast differentiation factors. *Biosci*  
768 *Rep.* 2019;39 3 doi:10.1042/BSR20182152.
- 769 50. Song Y, Altarejos J, Goodarzi MO, Inoue H, Guo X, Berdeaux R, et al. CRT3  
770 links catecholamine signalling to energy balance. *Nature.* 2010;468 7326:933-  
771 9. doi:10.1038/nature09564.
- 772 51. Liu J, Nong Q, Wang J, Chen W, Xu Z, You W, et al. Breed difference and  
773 regulatory role of CRT3 in porcine intramuscular adipocyte. *Anim Genet.*  
774 2020; doi:10.1111/age.12945.
- 775 52. Hoover KC. Smell with inspiration: the evolutionary significance of olfaction.  
776 *Am J Phys Anthropol.* 2010;143 Suppl 51:63-74. doi:10.1002/ajpa.21441.
- 777 53. Carneiro M, Rubin CJ, Di Palma F, Albert FW, Alfoldi J, Martinez Barrio A, et  
778 al. Rabbit genome analysis reveals a polygenic basis for phenotypic change  
779 during domestication. *Science.* 2014;345 6200:1074-9.  
780 doi:10.1126/science.1253714.
- 781 54. Tsao D, Thomsen HK, Chou J, Stratton J, Hagen M, Loo C, et al. TrkB agonists  
782 ameliorate obesity and associated metabolic conditions in mice. *Endocrinology.*  
783 2008;149 3:1038-48. doi:10.1210/en.2007-1166.

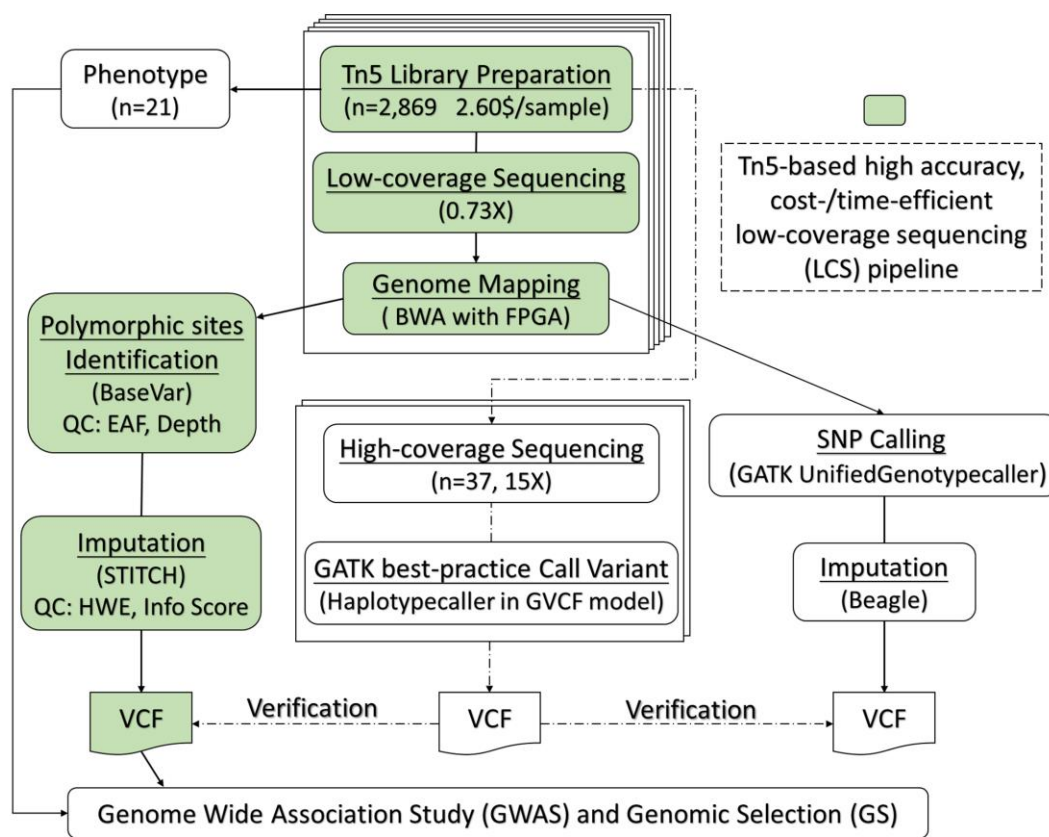
- 784 55. Lebrun B, Bariohay B, Moyse E and Jean A. Brain-derived neurotrophic factor  
785 (BDNF) and food intake regulation: a minireview. *Auton Neurosci.* 2006;126-  
786 127:30-8. doi:10.1016/j.autneu.2006.02.027.
- 787 56. Picelli S, Bjorklund AK, Reinius B, Sagasser S, Winberg G and Sandberg R.  
788 Tn5 transposase and tagmentation procedures for massively scaled sequencing  
789 projects. *Genome Res.* 2014;24 12:2033-40. doi:10.1101/gr.177881.114.
- 790 57. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernysky A, et  
791 al. The Genome Analysis Toolkit: a MapReduce framework for analyzing next-  
792 generation DNA sequencing data. *Genome Res.* 2010;20 9:1297-303.  
793 doi:10.1101/gr.107524.110.
- 794 58. Cingolani P, Platts A, Wang le L, Coon M, Nguyen T, Wang L, et al. A program  
795 for annotating and predicting the effects of single nucleotide polymorphisms,  
796 SnpEff: SNPs in the genome of *Drosophila melanogaster* strain w1118; iso-2;  
797 iso-3. *Fly (Austin).* 2012;6 2:80-92. doi:10.4161/fly.19695.
- 798 59. Yang J, Lee SH, Goddard ME and Visscher PM. GCTA: a tool for genome-wide  
799 complex trait analysis. *Am J Hum Genet.* 2011;88 1:76-82.  
800 doi:10.1016/j.ajhg.2010.11.011.
- 801 60. Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, et al. The  
802 variant call format and VCFtools. *Bioinformatics.* 2011;27 15:2156-8.  
803 doi:10.1093/bioinformatics/btr330.
- 804 61. Tajima F. Statistical method for testing the neutral mutation hypothesis by DNA  
805 polymorphism. *Genetics.* 1989;123 3:585-95.
- 806 62. Strimmer K. fdrtool: a versatile R package for estimating local and tail area-  
807 based false discovery rates. *Bioinformatics.* 2008;24 12:1461-2.  
808 doi:10.1093/bioinformatics/btn209.
- 809 63. Wang C, Prakapenka D, Wang S, Pulugurta S, Runesha HB and Da Y.  
810 GVCBLUP: a computer package for genomic prediction and variance  
811 component estimation of additive and dominance effects. *BMC Bioinformatics.*  
812 2014;15:270. doi:10.1186/1471-2105-15-270.



813 64. Capriotti E, Calabrese R and Casadio R. Predicting the insurgence of human  
814 genetic diseases associated to single point protein mutations with support vector  
815 machines and evolutionary information. *Bioinformatics*. 2006;22 22:2729-34.  
816 doi:10.1093/bioinformatics/btl423.

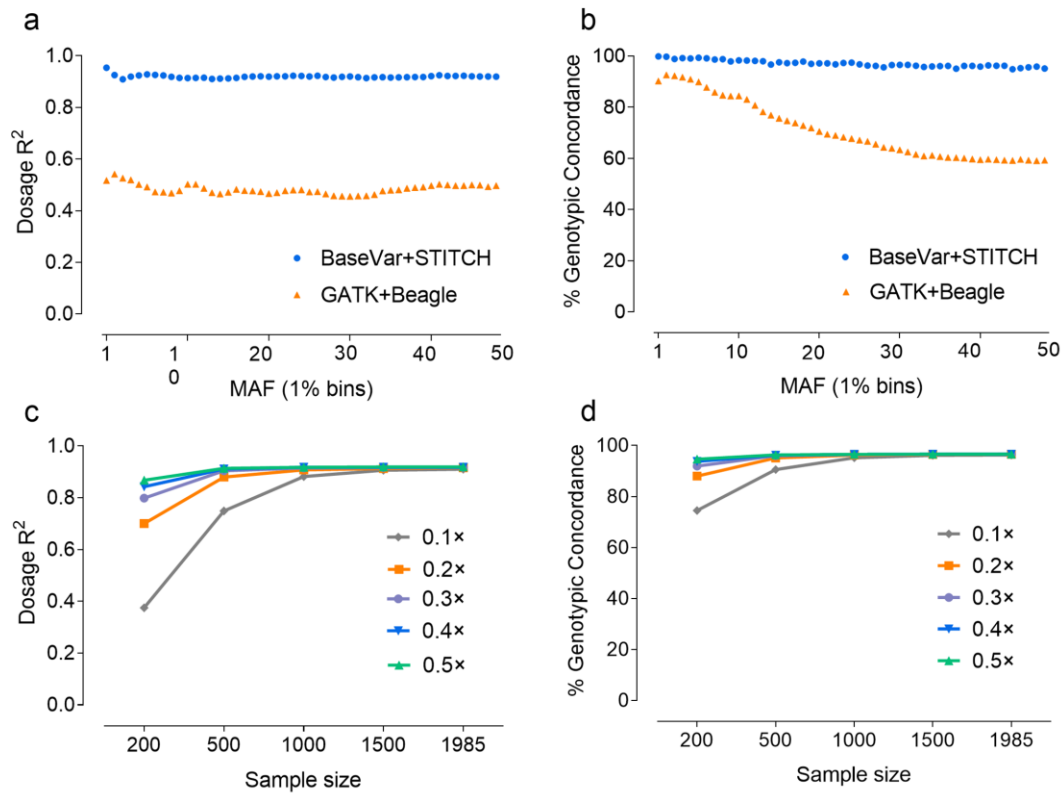
817

## Figure legends



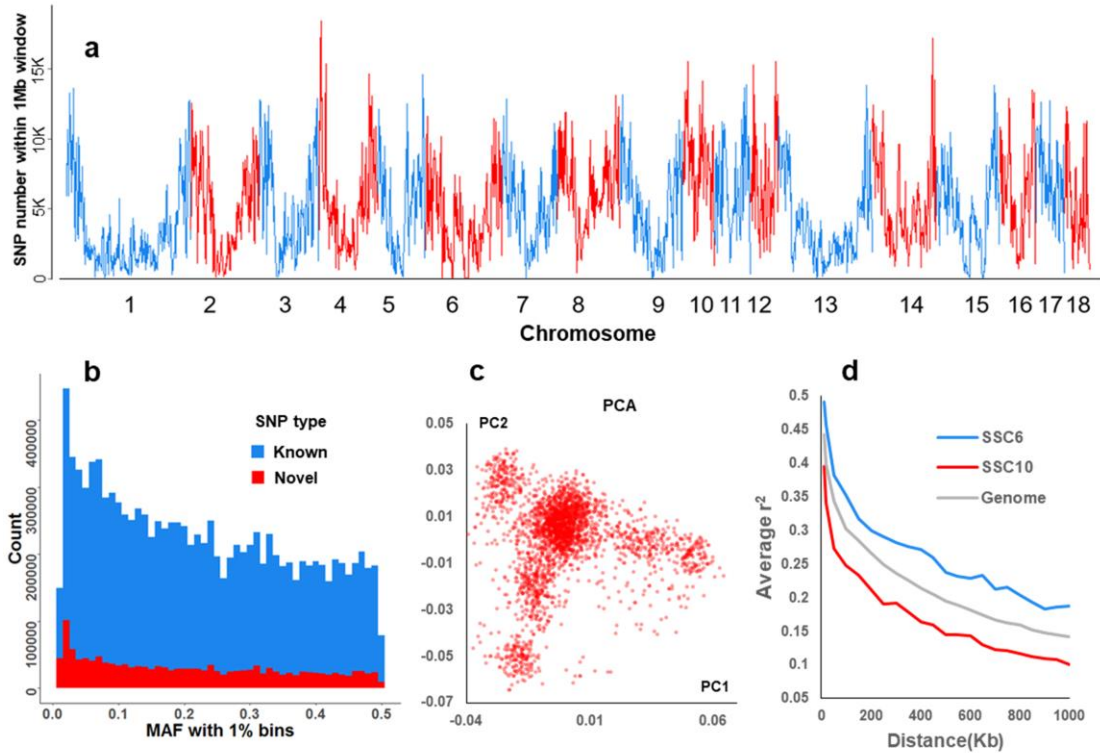
**Figure 1 The low coverage sequencing (LCS) study design**

The flow chart summarizes the steps used to identify and impute polymorphic sites, where the green block (left) represents the highly accurate pipeline used for the Tn5-based LCS analysis (BaseVar-STITCH). We also generated SNP results using the GATK-Beagle pipeline (right) and compared them with those found with the BaseVar-STITCH method. The data generated from the high-coverage sequencing analyses (middle) were used to assess the accuracy of the above results. The BaseVar-STITCH pipeline was used in the further GWAS presented in this study.



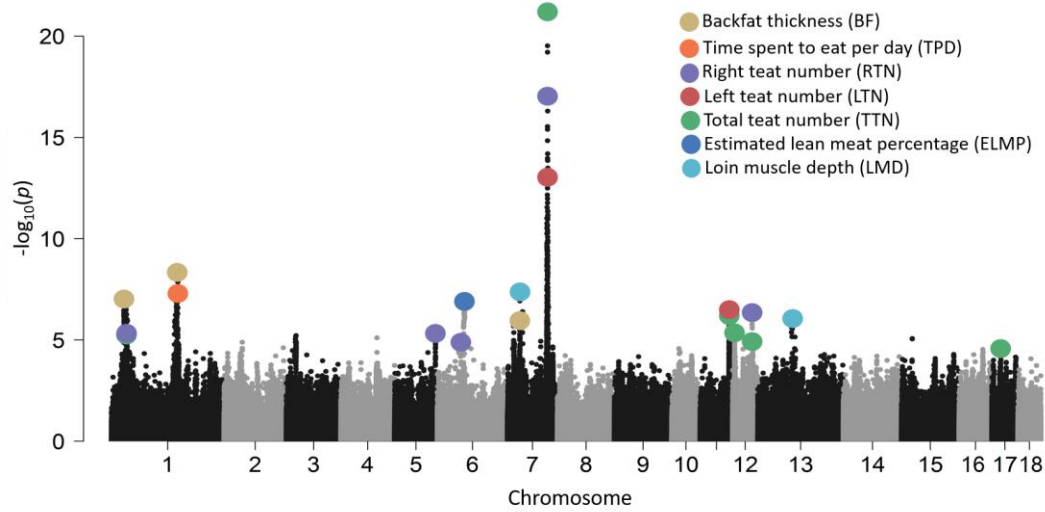
**Figure 2 Performance of BaseVar-STITCH on different minor allele frequencies (MAFs) and sample sizes**

The validation dataset is the high coverage sequencing results of 37 individuals genotyped by GATK best practices (HaplotypeCaller model). **(a)** and **(b)** show a comparison of the dosage  $R^2$  and genotypic concordance values (%) between the BaseVar-STITCH for low-coverage sequencing (LCS) (blue) and the GATK-Beagle (orange) pipelines, and **(c)** and **(d)** show the comparison of the dosage  $R^2$  and genotypic concordance values (%) among different sequencing depths.



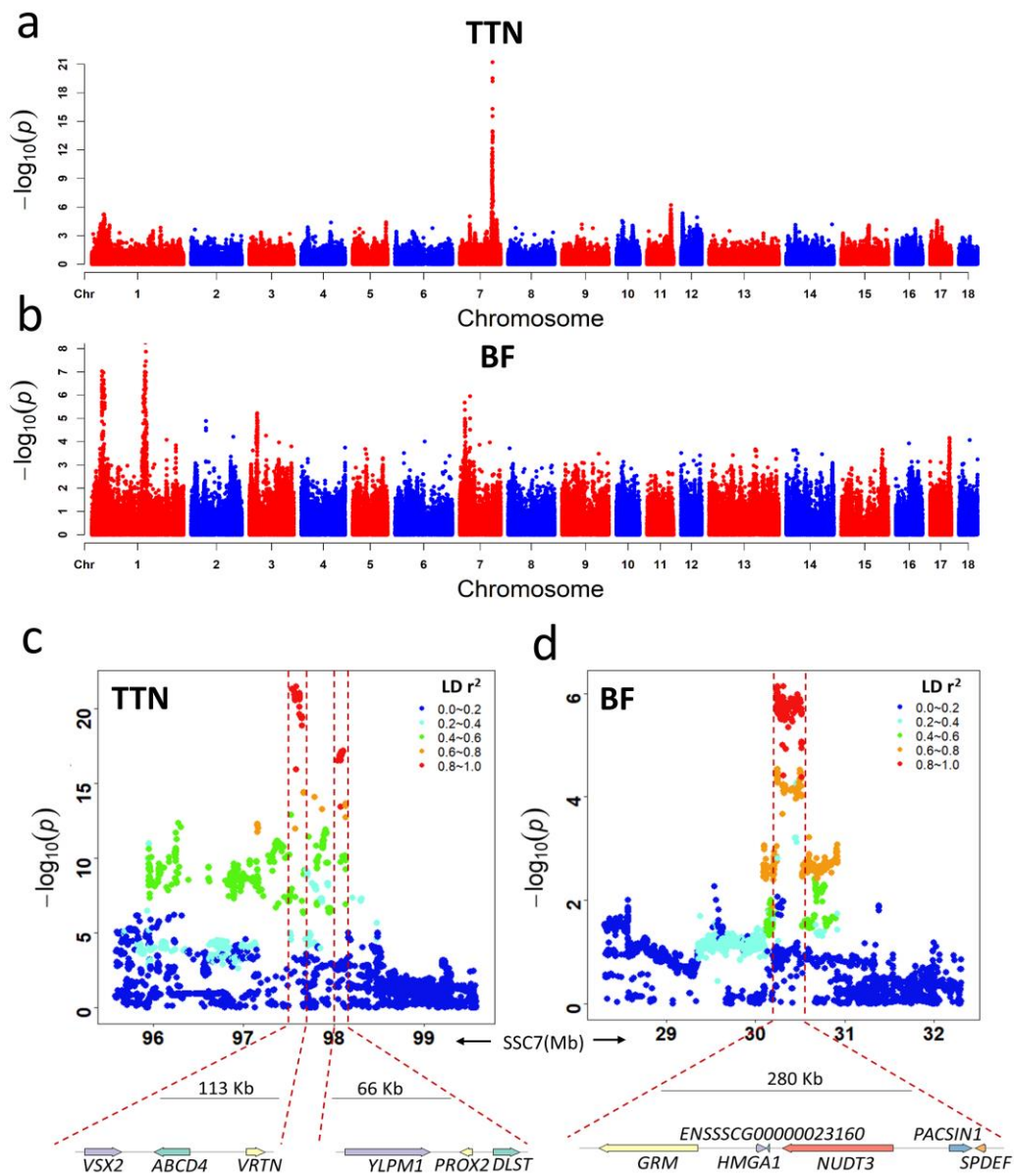
**Figure 3 Genetic diversity of the Duroc population**

(a) The distribution of SNPs in 1 Mb windows across the genome. (b) Histogram of allele counts by each 1% MAF bin. Novel (red) and known SNP sets (blue) were defined by comparing them to the pig dbSNP database. (c) Principal component 1 and 2 distribution in the Duroc population. (d) The extent of linkage disequilibrium (LD), in which the LD on chromosomes 6 (SSC6) and 10 (SSC10) represent the highest and lowest levels across the whole genome, respectively.



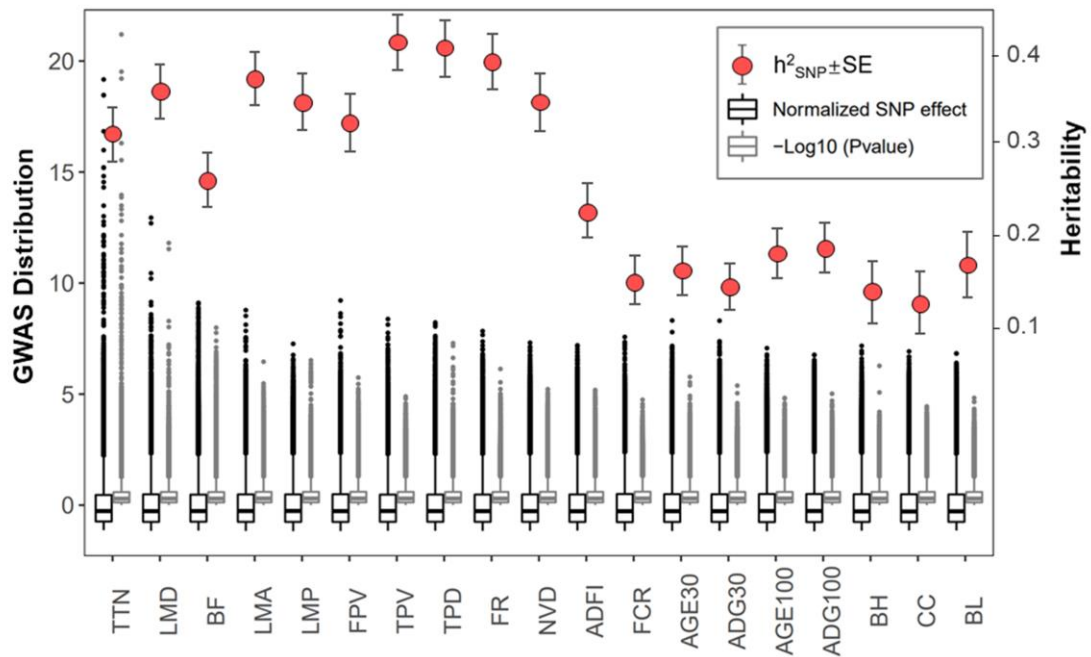
**Figure 4 Summary Manhattan plot of seven phenotypes with significant SNPs**

Genome-wide representation of all quantitative trait loci (QTLs) identified in this study. Light and dark grey dots show associations from the seven measures where at least one QTL was detected at the tagging SNP positions ( $n = 258,662$ ). The most significant SNP positions at each QTL are marked with a color dot.



**Figure 5** Manhattan plots and fine-mapping of the total test number (TTN) and back fat thickness (BF)

(a) and (b) Depict the TTN and BF association signals on the whole genome. (c) Fine-mapping of the TTN using the entire set of SNPs, in which two isolated regions on chromosome 7 with lengths of 113 and 66 Kb were detected as QTLs. (d) Fine-mapping of BF using the entire set of SNPs. A narrow QTL with a length of 280 Kb was detected on chromosome 7. The association genes within QTLs are displayed below.



**Figure 6 Heritability and SNP significance and normalized effect of 21 traits**

The SNP effect was estimated and normalized and is displayed in the black boxplot. The gray boxplot represents the distribution of  $-\log_{10} P$  values of all SNPs. Heritability estimates are represented by red dots, and black lines represent standard deviations.

**Table 1. QTLs mapping and contribution to heritability**

Phenotype	Number	Mean $\pm$ standard deviation	Significant threshold <sup>a</sup>	QTL number	Variance explained(%) <sup>b</sup>	Gene number <sup>c</sup>
Total teat number (TTN)	2797	10.73 $\pm$ 1.07	4.55	6	8.86	52
Left teat number (LTN)	2797	5.35 $\pm$ 0.66	4.81	2	3.16	14
Right teat number (RTN)	2797	5.38 $\pm$ 0.64	4.79	5	6.03	56
Back fat thickness at 100 Kg (BF, mm)	2796	10.99 $\pm$ 2.66	4.67	4	2.40	55
Loin muscle depth at 100 Kg (LMD, mm)	2796	46.15 $\pm$ 3.93	5.36	2	1.27	15
Loin muscle area at 100 Kg (LMA, mm <sup>2</sup> )	2795	36.25 $\pm$ 3.60	-	0	0	0
Lean meat percentage at 100 Kg (LMP, %)	2795	54.02 $\pm$ 1.58	5.50	1	1.19	48
Time spent to eat per day (TPD, min)	2602	63.02 $\pm$ 9.85	6.10	1	1.08	28
Average daily feed intake (ADFI, Kg)	2602	2.00 $\pm$ 0.20	-	0	0	0
Number of visits to feeder per day (NVD)	2602	7.30 $\pm$ 1.83	-	0	0	0
Time spent to eat per visit (TPV, min)	2602	10.06 $\pm$ 2.79	-	0	0	0
Feed intake rate (FR, g/min)	2602	32.37 $\pm$ 5.19	-	0	0	0
Feed intake per visit (FPV, Kg)	2602	290.6 $\pm$ 75.87	-	0	0	0
Feed conversion rate (FCR)	2691	2.19 $\pm$ 0.19	-	0	0	0
Average daily gain (0-30 Kg) (ADG30, g)	2795	354.8 $\pm$ 38.72	-	0	0	0
Age to 30 kg live weight (AGE30, day)	2796	80.49 $\pm$ 8.57	-	0	0	0
Average daily gain (30-100 Kg) (ADG100, g)	2795	633.8 $\pm$ 37.12	-	0	0	0
Age to 100 kg live weight (AGE100, day)	2796	155.5 $\pm$ 9.20	-	0	0	0
Body length (BL, cm)	1844	117.60 $\pm$ 2.91	-	0	0	0
Body height (BH, cm)	1844	62.19 $\pm$ 1.55	-	0	0	0
Circumference of cannon bone (CC, cm)	1844	17.81 $\pm$ 0.54	-	0	0	0

a.  $-\text{Log}_{10}(p)$  value when  $\text{FDR} < 0.05$ ; b. total phenotypic variance explained by QTLs; c. Total gene number included in QTLs.



## **Additional Files**

### **Supplementary Figure 1 Dosage $R^2$ and cost time (minute) among different K values**

Accuracy and cost time of genotyping from  $K = 5$  to  $K = 25$ , where the blue and black lines represent the dosage  $R^2$  and cost time (minute) respectively.

### **Supplementary Figure 2 Purifying selection regions in the whole genome**

Purifying selection signals were detected on SSC2, SSC3, SSC6, SSC7, SSC9 and SSC15, where blue and red lines represent  $-\text{Log}_{10} \text{Pi}$  and Tajima's  $D$  respectively, and the grey regions depict the purifying selected regions.

### **Supplementary Figure 3 Phenotypic distribution of 21 traits**

### **Supplementary Figure 4 Manhattan plots of phenotypes with no significant SNPs**

Manhattan plots of ADFT, NVD, TPV, FPV, FR, FCR, BH, BL, CC, ADG100, AGE100, ADG30, AGE30 and LMA, where no significant SNPs were detected in these traits.

### **Supplementary Figure 5 QQ plot of 21 phenotypes**

### **Supplementary Figure 6 Summary plots of fine mapping**

### **Supplementary Figure 7 Distribution of top 100 SNPs based on $P$ value using GWAS analysis**

### **Supplementary Figure 8 GO and KEGG enrichment of genes identified to be associated with feeding behavior traits**

### **Supplementary Figure 9 Neurotrophin signaling pathway enrichment**

The red tangles represent detected pathways in this study, which including Bcl-2, NT3, TrkB and p75NTR.

### **Supplementary Table S1 LC data set**

### **Supplementary Table S2 Resequencing Duroc samples list**

### **Supplementary Table S3 Number and density of SNPs imputed by STITCH and Tag SNP**

**Supplementary Table S4 GO enrichment of genes located in the selected regions**

**Supplementary Table S5 Summary of detected QTLs**

**Supplementary Table S6 Summary table of markers identified significantly associated with ADG, AGE or FCR in previous studies**

**Supplementary Table S7 GO enrichment of genes located in the selected regions**

**Supplementary Table S8 KEGG enrichment of genes located in the selected regions**

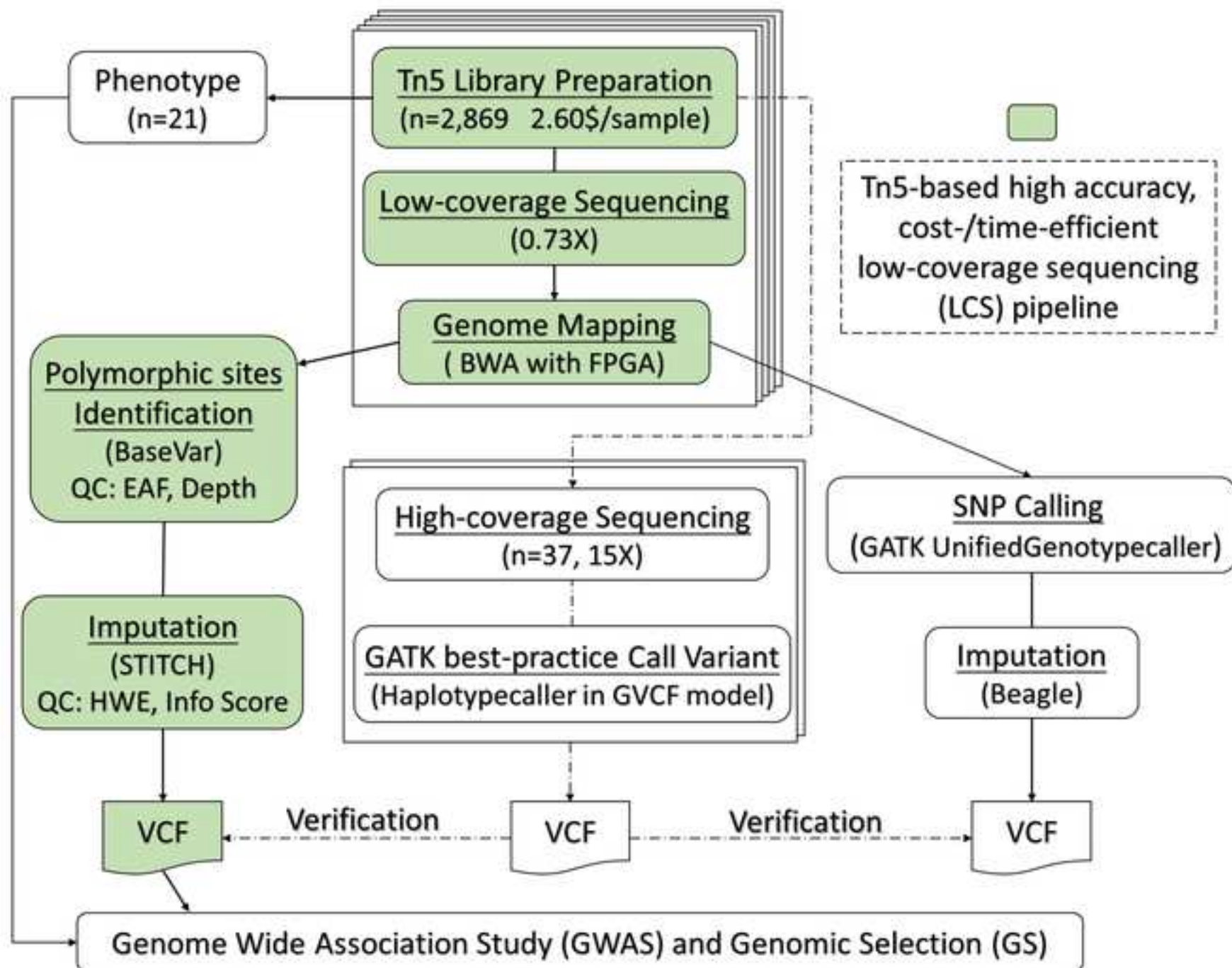
**Supplementary Table S9 Missense SNPs in the narrowed QTL region of TN**

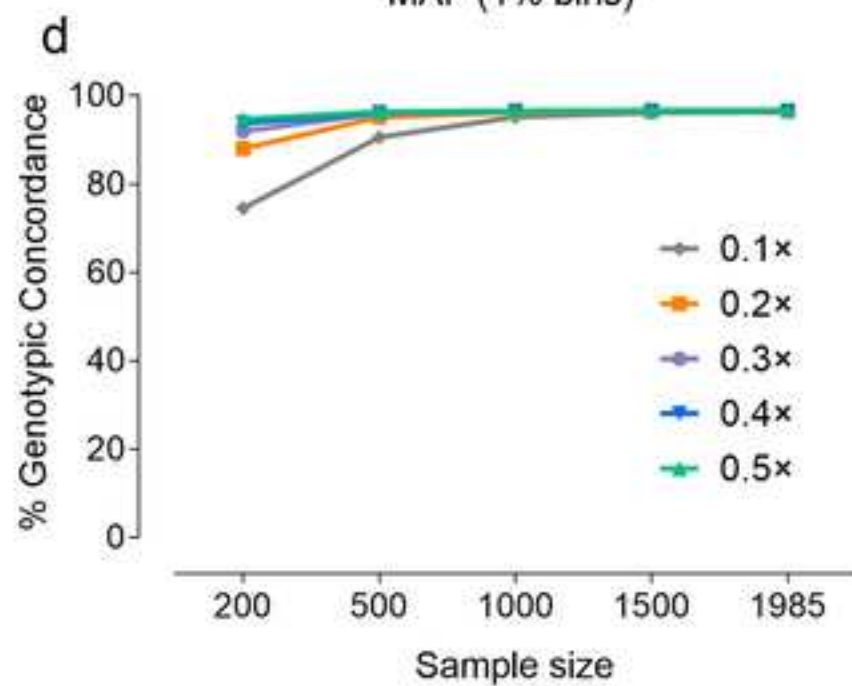
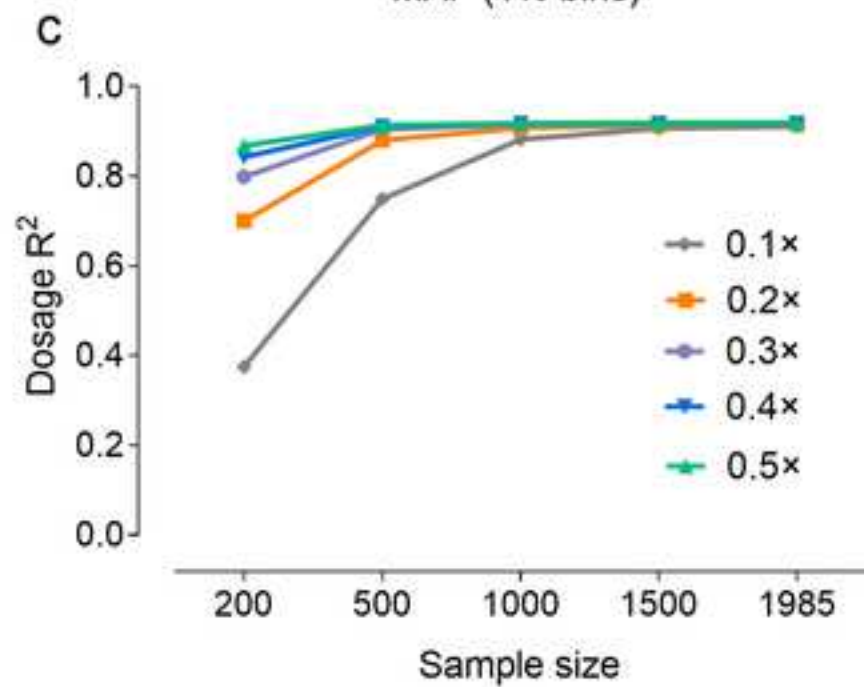
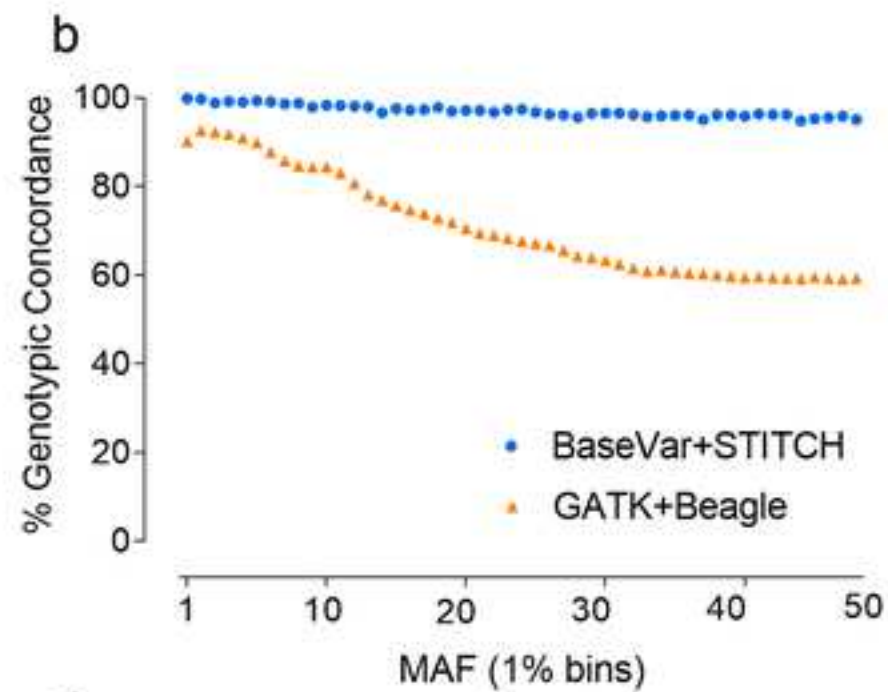
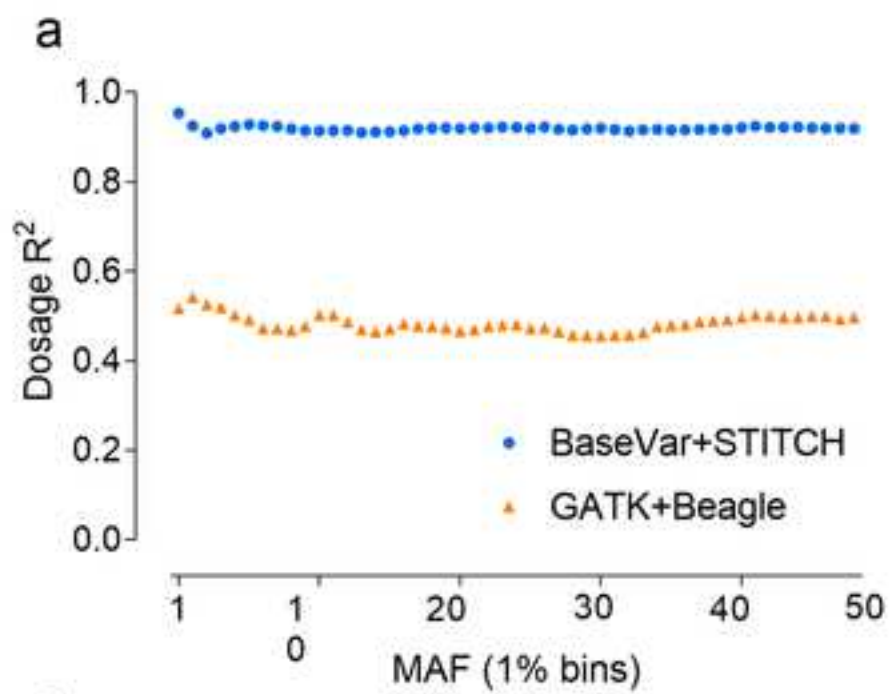
**Supplementary Table S10 Gathered information of candidate genes**

**Table 1 QTLs mapping and contribution to heritability**

Phenotype	Number	Mean $\pm$ standard deviation	Significant threshold <sup>a</sup>	QTL number	Variance explained(%) <sup>b</sup>	Gene number <sup>c</sup>
Total teat number (TTN)	2797	10.73 $\pm$ 1.07	4.55	6	8.86	52
Left teat number (LTN)	2797	5.35 $\pm$ 0.66	4.81	2	3.16	14
Right teat number (RTN)	2797	5.38 $\pm$ 0.64	4.79	5	6.03	56
Back fat thickness at 100 Kg (BF, mm)	2796	10.99 $\pm$ 2.66	4.67	4	2.40	55
Loin muscle depth at 100 Kg (LMD, mm)	2796	46.15 $\pm$ 3.93	5.36	2	1.27	15
Loin muscle area at 100 Kg (LMA, mm <sup>2</sup> )	2795	36.25 $\pm$ 3.60	-	0	0	0
Lean meat percentage at 100 Kg (LMP, %)	2795	54.02 $\pm$ 1.58	5.50	1	1.19	48
Time spent to eat per day (TPD, min)	2602	63.02 $\pm$ 9.85	6.10	1	1.08	28
Average daily feed intake (ADFI, Kg)	2602	2.00 $\pm$ 0.20	-	0	0	0
Number of visits to feeder per day (NVD)	2602	7.30 $\pm$ 1.83	-	0	0	0
Time spent to eat per visit (TPV, min)	2602	10.06 $\pm$ 2.79	-	0	0	0
Feed intake rate (FR, g/min)	2602	32.37 $\pm$ 5.19	-	0	0	0
Feed intake per visit (FPV, Kg)	2602	290.6 $\pm$ 75.87	-	0	0	0
Feed conversion rate (FCR)	2691	2.19 $\pm$ 0.19	-	0	0	0
Average daily gain (0-30 Kg) (ADG30, g)	2795	354.8 $\pm$ 38.72	-	0	0	0
Age to 30 kg live weight (AGE30, day)	2796	80.49 $\pm$ 8.57	-	0	0	0
Average daily gain (30-100 Kg) (ADG100, g)	2795	633.8 $\pm$ 37.12	-	0	0	0
Age to 100 kg live weight (AGE100, day)	2796	155.5 $\pm$ 9.20	-	0	0	0
Body length (BL, cm)	1844	117.60 $\pm$ 2.91	-	0	0	0
Body height (BH, cm)	1844	62.19 $\pm$ 1.55	-	0	0	0
Circumference of cannon bone (CC, cm)	1844	17.81 $\pm$ 0.54	-	0	0	0

*Note:* a.  $-\log_{10} P$  value when FDR < 0.05; b. total phenotypic variance explained by QTLs; c. Total gene number included in QTLs.





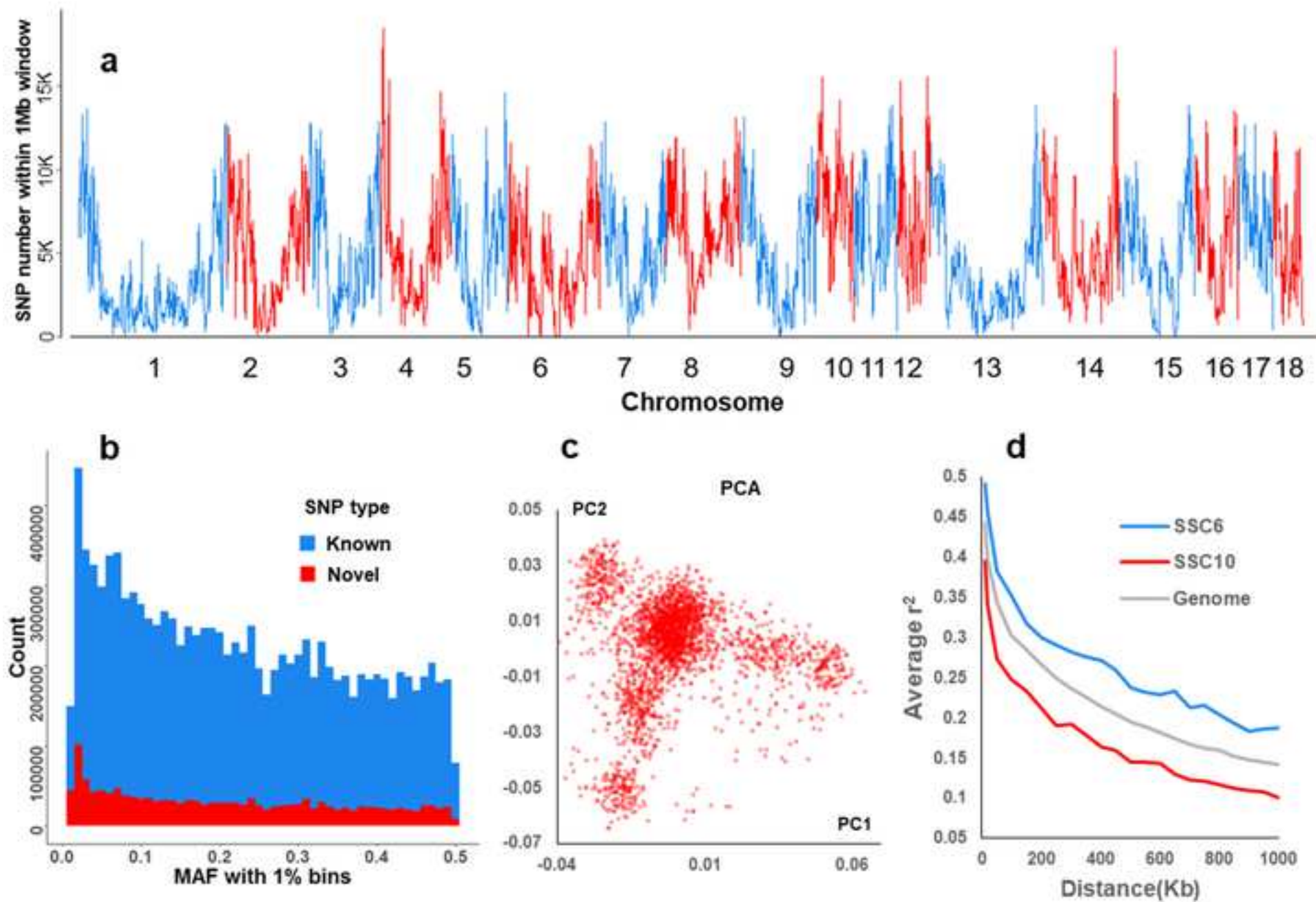
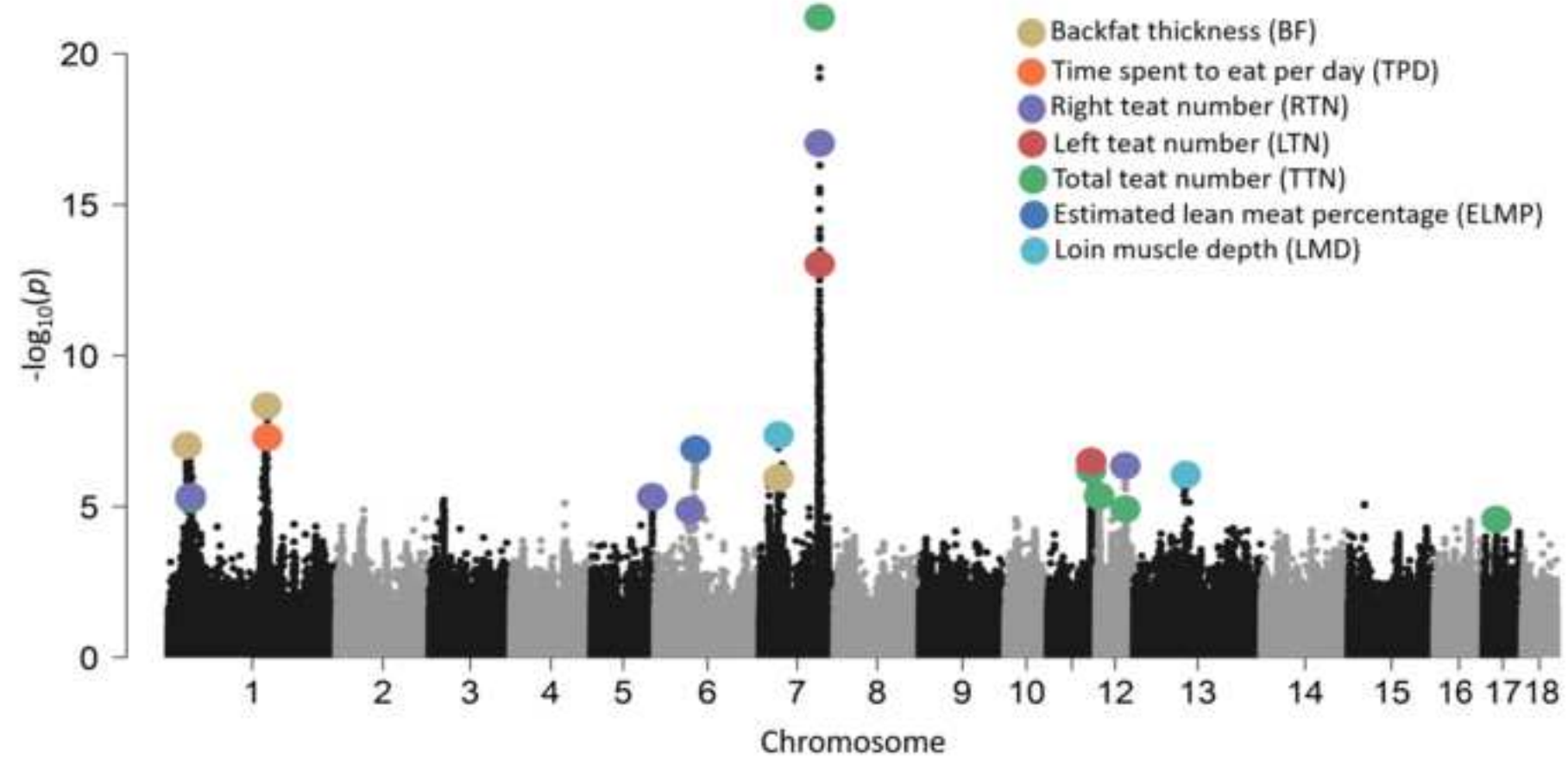
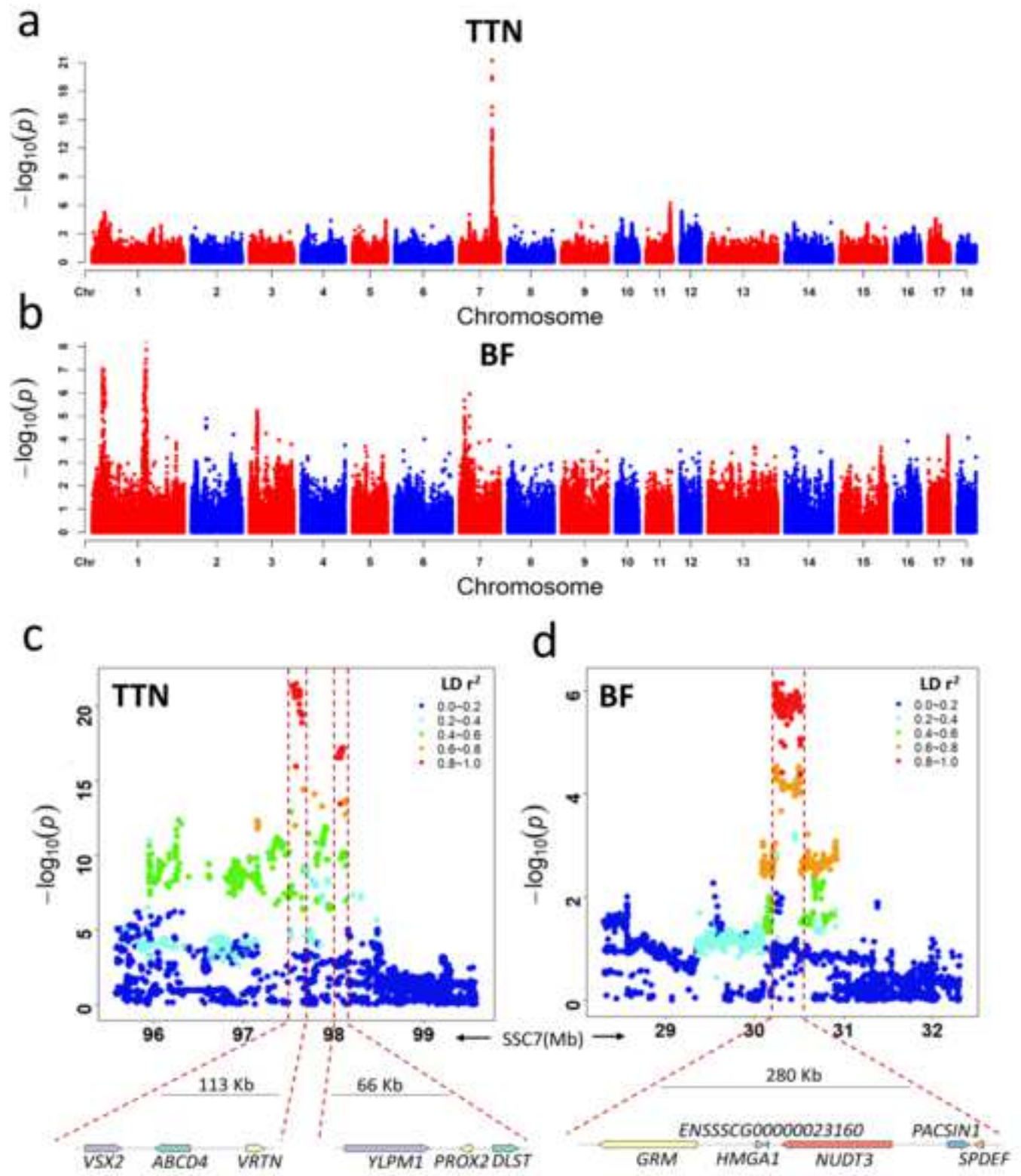


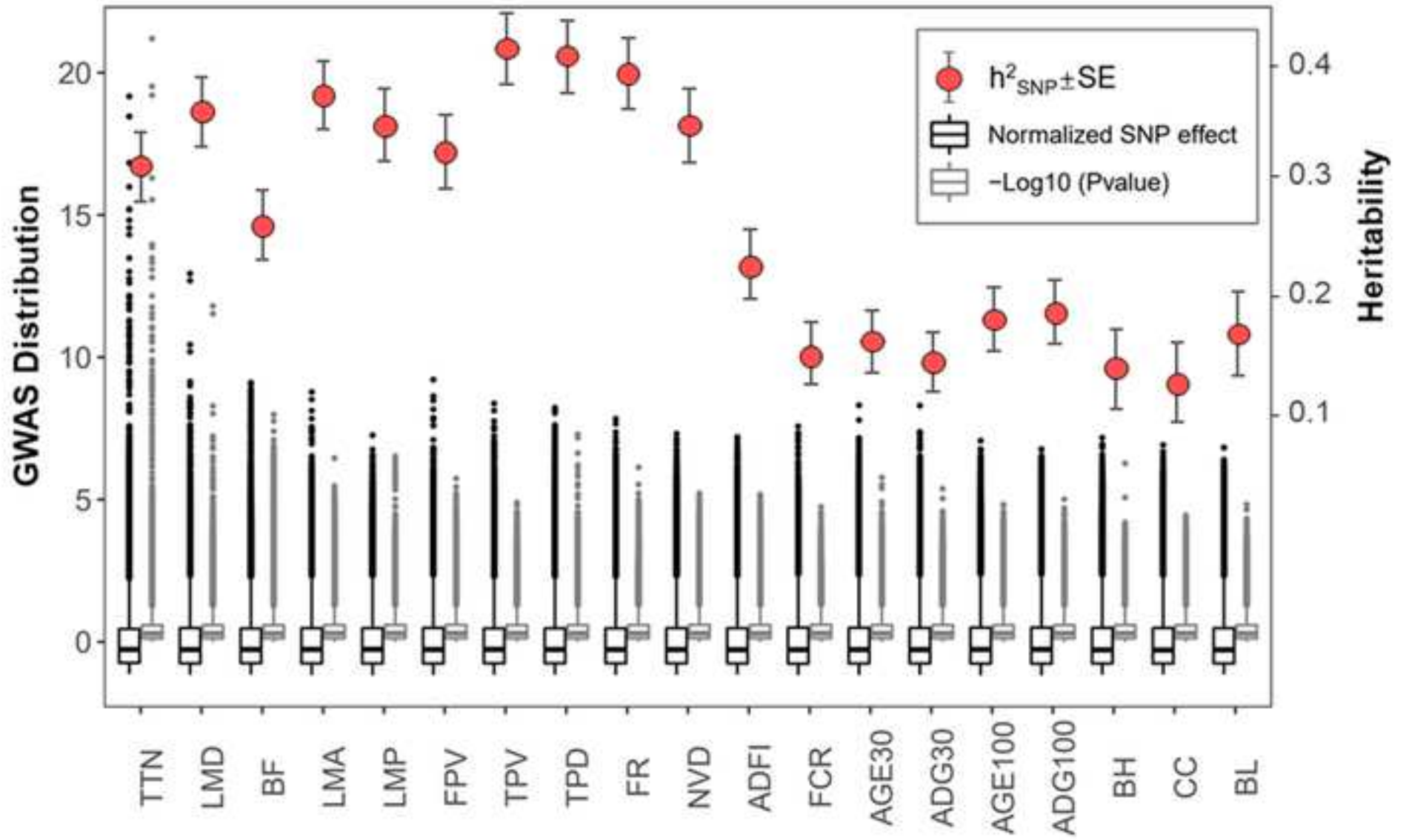


Figure 4











Click here to access/download  
**Supplementary Material**  
Supplementary Figure 1.tif













Click here to access/download  
**Supplementary Material**  
Supplementary Figure 6.tif











Click here to access/download  
**Supplementary Material**  
Supplementary Figure 9.png





Click here to access/download  
**Supplementary Material**  
Supplementary Table S1.xlsx



Click here to access/download  
**Supplementary Material**  
Supplementary Table S2.xlsx





Click here to access/download  
**Supplementary Material**  
Supplementary Table S3.xlsx





Click here to access/download  
**Supplementary Material**  
Supplementary Table S4.xlsx





Click here to access/download  
**Supplementary Material**  
Supplementary Table S5.xlsx



Click here to access/download  
**Supplementary Material**  
Supplementary Table S6.xlsx







Click here to access/download  
**Supplementary Material**  
Supplementary Table S7.xlsx





Click here to access/download  
**Supplementary Material**  
Supplementary Table S8.xlsx

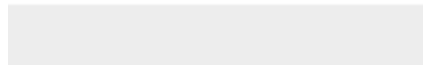


Click here to access/download  
**Supplementary Material**  
Supplementary Table S9.xlsx





Click here to access/download  
**Supplementary Material**  
SupplementaryTable S10.xlsx



## Cover Letter

December 01, 2020

Dear Editor,

We would like to resubmit to *GigaScience* the modified manuscript entitled “Accelerated Deciphering of the Genetic Architecture of Agricultural Economic Traits in Pigs Using the Low Coverage Whole-genome Sequencing Strategy”. We believe that this manuscript will make it interesting to general readers of your journal.

Domestication not only modified the economic important traits but also left a genetic signature that affects both the population diversity and genomic structure of domesticated farm animals. Fully elucidating the phenotypic diversity, revealing the genetic structure of the breeding population is the basis for precision breeding. Large-scale WGS and GWAS strategies had enable us to gain different perspectives which was not possible before. However, high depth sequencing in large cohorts is still prohibitively expensive, to develop a massively parallel low coverage sequencing method has become imperative.

Here, we report a Tn5-based, highly accurate, cost and time-efficient, low coverage sequencing (LCS) approach to perform sequencing on 2,869 Duroc boars at an average depth of 0.73×, which identify 11.3 M SNPs throughout the genome. Base on the whole genome sequencing strategy, the high-resolution genome-wide association study (GWAS) detected 14 candidate quantitative trait loci (QTLs) in 7 of 21 important traits and provided a lot of worth points for further investigation. We also showed that the artificial selection alters genomes that affect important growth traits. Moreover, we explored the different traits with varies genetic architecture in depth, providing guidance for subsequent genetic improvement by genomic selection. The LCS strategy, together with the unprecedented capacity of NGS allows the cost-effective and large-scale genome analysis with industrial-scale efficiency, and we are also confident that it will be a universal strategy to meet the needs for the genomic study and breeding of both animals and plants.

All of the sequencing raw data in this study have been deposited into NCBI with accession number PRJNA681437, the variance data as VCF file will be available via GIGADB. The data will be shared publicly without restrictions in case of acceptance.

We confirm that this manuscript has not been published elsewhere and is not under consideration by another journal, and all authors declare that they have no competing interests.

Thank you for your consideration. We look forward to hearing from you at your earliest convenience.

Xiaoxiang Hu

The State Key Laboratory for Agro-biotechnology, China Agricultural University,  
Beijing, PRC,100193

Email: [huxx@cau.edu.cn](mailto:huxx@cau.edu.cn)