

# Optimizing Genomic Selection in Dezhou Donkey Using Low Coverage Whole Genome Sequencing

**Changheng Zhao**

Shandong Agricultural University

**Jun Teng**

Shandong Agricultural University

**Xinhao Zhang**

Shandong Agricultural University

**Dan Wang**

Shandong Agricultural University

**Xinyi Zhang**

Shandong Agricultural University

**Shiyin Li**

Shandong Agricultural University

**Haijing Li**

Dong-E E-Jiao Co.

**Xin Jiang**

Shandong Agricultural University

**Chao Ning**

Shandong Agricultural University

**Qin Zhang** (✉ [qzhang@cau.edu.cn](mailto:qzhang@cau.edu.cn))

Shandong Agricultural University

---

## Research

**Keywords:** Dezhou donkey, Low coverage whole genome sequencing, Genotype imputation, Genomic selection

**Posted Date:** June 17th, 2021

**DOI:** <https://doi.org/10.21203/rs.3.rs-607740/v1>

**License:** © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

1 **Optimizing Genomic Selection in Dezhou Donkey Using low**  
2 **coverage whole genome sequencing**

3 Changheng Zhao<sup>1</sup>, Jun Teng<sup>1</sup>, Xinhao Zhang<sup>1,2</sup>, Dan Wang<sup>1</sup>, Xinyi Zhang<sup>1</sup>, Shiyin Li<sup>1</sup>,  
4 Xin Jiang<sup>1</sup>, Haijing Li<sup>2</sup>, Chao Ning<sup>1\*</sup>, Qin Zhang<sup>1\*</sup>

5 1. Shandong Provincial Key Laboratory of Animal Biotechnology and Disease Control  
6 and Prevention, College of Animal Science and Veterinary Medicine, Shandong  
7 Agricultural University, Tai'an, 271018, China

8 2. National Engineering Research Center for Gelatin-based TCM, Dong-E E-Jiao Co.,  
9 Ltd, 78 E-Jiao Street, Donge County, 252201, Shandong Province, China

10  
11  
12  
13 \* Corresponding author:

14 Chao Ning

15 ningchao@sdau.edu.cn

16  
17 Qin Zhang

18 qzhang@sdau.edu.cn  
19  
20  
21  
22

23 **Abstract**

24 **Background:** Low coverage whole genome sequencing is a low-cost genotyping  
25 technology. Combining with genotype imputation approaches, it is likely to become a  
26 critical component of cost-efficient genomic selection programs in agricultural  
27 livestock. Here, we used the low-coverage sequence data of 617 Dezhou donkeys to  
28 investigate the performance of genotype imputation for low coverage whole genome  
29 sequence data and genomic selection based on the imputed genotype data. The specific  
30 aims were: (i) to measure the accuracy of genotype imputation under different  
31 sequencing depths, sample sizes, MAFs, and imputation pipelines; and (ii) to assess the  
32 accuracy of genomic selection under different marker densities derived from the  
33 imputed sequence data, different strategies for constructing the genomic relationship  
34 matrixes, and single- vs multi-trait models.

35 **Results:** We found that a high imputation accuracy ( $> 0.95$ ) can be achieved for  
36 sequence data with sequencing depth as low as 1x and the number of sequenced  
37 individuals equal to 400. For genomic selection, the best performance was obtained by  
38 using a marker density of 410K and a **G** matrix constructed using marker dosage  
39 information. Multi-trait GBLUP performed better than single-trait GBLUP.

40 **Conclusions:** Our study demonstrates that low coverage whole genome sequencing  
41 would be a cost-effective method for genomic selection in Dezhou Donkey.

42 **Keywords:** Dezhou donkey, Low coverage whole genome sequencing, Genotype  
43 imputation, Genomic selection

44

## 45 **Background**

46 Dezhou Donkey, originated from Dezhou area, Shandong Province, China, is one  
47 of major donkey breeds in China. It is famous for its large body size (thus good meat  
48 production ability) and excellent skin quality (for producing donkey-hide gelatin). It  
49 has been introduced as breeding stock into many provinces and cities, and has also  
50 brought considerable economic benefits to farmers. Therefore, Dezhou Donkey plays  
51 an important role in the donkey industry in China. However, selective breeding in the  
52 sense of modern animal breeding theory has long been ignored for any donkey breeds  
53 in China. In recent years, along with the increasing of the importance of the donkey  
54 industry in livestock agriculture in China, donkey breeding is gradually becoming an  
55 important issue in donkey production and some breeding work are carrying out in  
56 Dezhou Donkey population.

57 Starting with the pioneered work of Meuwissen et al. [1], genomic selection (GS)  
58 has been widely used in selective breeding in almost all major farm animal species, and  
59 has brought great increasement of genetic progresses and economic benefit for many  
60 animal breeding industries [2-4] . Typically, GS is carried out using a high (or medium)  
61 density marker (SNP) array. Many commercial SNP arrays have been developed for  
62 almost all major farm animal species. However, there is no such array for donkey, which  
63 inhibits the application of GS in donkey.

64 Recently, along with the rapid development of next generation sequencing  
65 technology and reduction of sequencing cost, GS using genotypes revealed by whole  
66 genome sequencing (WGS, instead of SNP array) has drawn interests of animal GS

67 community with the motivation of further improving the selection accuracy, better  
68 application of GS across breeds/populations, and better persistence of accuracy across  
69 generations [5-6]. To capture all variants in the genome, a sequencing depth of about  
70 10x is generally required. However, at present 10x sequencing is still too expensive for  
71 large scale GS application. An alternative is to perform low coverage whole genome  
72 sequencing (lcWGS) at only about 1x, and then recovering the missing genotypes by  
73 imputation to ensure that all individuals have genotypes for a shared set of variants.  
74 This approach has been used in human and some animal species for genome-wide  
75 association study and genomic selection/prediction and approved to be a feasible  
76 alternative to normal sequencing [7-10]. Since the cost of lcWGS can even be lower  
77 than that of SNP array, it is considered as a cost-effective genotyping approach for GS  
78 and GS based this approach was referred as GS 2.0 by Hickey (2013) [11].

79 A critical issue of lcWGS-based GS is the accuracy of imputation of missing  
80 genotypes, which is affected by several factors, such as sequencing depth, population  
81 size, minor allele frequency (MAF), and imputation method. A number of imputation  
82 methods for lcWGS data have been proposed [12-14]. However, most of these methods  
83 require a high-density haplotype reference panel, which are not available for most  
84 animal species. Davies et al. (2016)[12] proposed a method called STITCH for  
85 imputation based only on sequencing read data, without requiring a haplotype reference  
86 panel, which provides an opportunity of using lcWGS technology for species that lack  
87 a haplotype reference panel.

88 In this study, we evaluated the imputation accuracy of lcWGS data with respect to

89 different sequencing depths, population sizes, MAFs, and imputation pipelines using  
90 617 Dezhou Donkey animals which were sequenced with an average depth of 3.5x. We  
91 then used the imputed genotypes to investigate the performance of genomic selection  
92 for birth weight and weaning weight in the Dezhou Donkey population under different  
93 marker densities, strategies for constructing the genomic relationship matrices, and  
94 single- vs two-trait models.

95

## 96 **Materials and Methods**

### 97 *Animals*

98 Blood samples from 617 Dezhou Donkey animals were collected from a donkey farm  
99 in Shandong Province. Total DNA was isolated using the QIAamp DNA Investigator  
100 kit (QIAGEN, Hilden, Germany) and following the manufacturer's instruction. DNA  
101 quality was evaluated by spectrophotometry and agarose gel electrophoresis.

102 All experimental chickens were maintained, and all the studies were carried out  
103 according to the guideline of the experimental animal management of Shandong  
104 Agricultural University (SDAUA-2018-018).

105

### 106 *Low coverage whole genome sequencing*

107 DNA templates were ultrasonically sheared using a Covaris E220 (Covaris,  
108 Woburn, MA, USA) to yield ~150 bp fragments, and then prepared for sequencing  
109 libraries following the workflow of the NEBNext Ultra DNA Library Preparation  
110 Protocol. Multiple Ampure Bead XP cleanups (Beckman Coulter, Brea, CA, USA) were

111 conducted to remove any adapter dimer that might have developed. The quality and  
112 concentration of libraries were determined on an Agilent Bioanalyzer 2100 (Agilent  
113 Technologies, Santa Clara, CA). The quality-controlled genomic library for each  
114 sample was PE150 sequenced using the Illumina NovaSeq 6000 sequencing system.  
115 The average sequencing coverage of the sample data was 3.5x.

116 Read quality was assessed using the FastQC software  
117 (<https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>) with focus on base  
118 quality scores and GC content, N content, and sequence duplication levels. All data  
119 reached a nucleotide length of longer than 50 bp and a Phred quality score of lower than  
120 30, which were aligned to the donkey reference genome [15] by BWA [16]. Samtools  
121 [17] was used to transfer formats, sort and index files.

122

### 123 *Pipelines for genotype imputation*

124 We compared two imputation pipelines, i.e., Bcftools + Beagle and BaseVar +  
125 STITCH. In the first pipeline, we called SNPs using Bcftools[18], performed quality  
126 control using PLINK [19] with the parameters of geno (>90%) and MAF (>1%), and  
127 then conducted genotype imputation using Beagle v4.1 [20]. In the second pipeline, we  
128 called SNPs using BaseVar [7] , filtered with EAF >= 0.01, and then imputed the  
129 missing genotypes (with probabilities) using STITCH. The resulted SNP data were  
130 filtered with an imputation info\_score > 0.4 and a Hardy-Weinberg Equilibrium (HWE)  
131  $p$ -value > 1e-6.

132

133 *Evaluation of imputation accuracy*

134 We evaluated the imputation accuracy using the sequencing data of 18 Dezhou  
135 Donkey animals, which were sequenced with an average sequencing coverage of 13.5x.  
136 Chromosomes 1, 19 and 30 were chosen to compare the imputed and typed genotypes  
137 in terms of genotype concordance measured as proportion of correctly imputed  
138 genotypes and genotype accuracy measured as squared Pearson correlation coefficient  
139 ( $r^2$ ) between imputed dosages and typed genotypes. To evaluate the imputation  
140 accuracy for different sequencing depth, we randomly sampled reads from the BAM  
141 files to generate sequence data with different sequencing depth (1x and 1.5x) using  
142 Picard (<https://broadinstitute.github.io/picard/>). The effects of sample size (number of  
143 low coverage sequenced individuals) and minor allele frequency (MAF) on the  
144 imputation accuracy were also tested.

145

146 *Genomic selection*

147 The imputation-based sequence data was used to investigate the performance of  
148 genomic prediction in Dezhou Donkey population. Two traits were considered, birth  
149 weight (BW) and weaning weight (WW). 594 animals with records on both traits and  
150 sequence data were included. The genomic breeding values were estimated using the  
151 genomic best linear unbiased prediction (GBLUP) [21] method under single-trait model  
152 as well as two-trait model.

153 Single-trait model:

154 
$$\mathbf{y} = \mathbf{Xb} + \mathbf{Za} + \mathbf{e}$$



155 where  $\mathbf{y}$  is the vector of observed phenotypes,  $\mathbf{b}$  is the vector of fixed effects, which  
 156 include sex effects and year-season effects,  $\mathbf{a}$  is the vector of genomic breeding values  
 157 with distribution of  $N(0, \mathbf{G}\sigma_a^2)$ , where  $\sigma_a^2$  is the additive genetic variance and  $\mathbf{G}$  is the  
 158 genomic relationship matrix,  $\mathbf{X}$  and  $\mathbf{Z}$  are the incidence matrices for  $\mathbf{b}$  and  $\mathbf{a}$ ,  
 159 respectively, and  $\mathbf{e}$  is the vector of random residuals with distribution of  $N(0, \mathbf{I}\sigma_e^2)$ .

160 Two-trait model:

$$161 \begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{X}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{X}_2 \end{bmatrix} \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{bmatrix} + \begin{bmatrix} \mathbf{Z}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{Z}_2 \end{bmatrix} \begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \end{bmatrix} + \begin{bmatrix} \mathbf{e}_1 \\ \mathbf{e}_2 \end{bmatrix}$$

162 where the meanings of the vectors and matrices are the same as those in the single-trait  
 163 model with the subscripts 1 and 2 referring trait 1 and trait 2, respectively. It was

164 assumed that  $\begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \end{bmatrix} \sim N(\mathbf{0}, \mathbf{G} \otimes \mathbf{M})$ , where  $\mathbf{M} = \begin{bmatrix} \sigma_{a_1}^2 & \sigma_{a_{12}} \\ \sigma_{a_{12}} & \sigma_{a_2}^2 \end{bmatrix}$  is the variance-covariance

165 matrix of the genomic breeding values of the two traits, and  $\begin{bmatrix} \mathbf{e}_1 \\ \mathbf{e}_2 \end{bmatrix} \sim N(\mathbf{0}, \mathbf{I} \otimes \mathbf{R})$ , where

166  $\mathbf{R} = \begin{bmatrix} \sigma_{e_1}^2 & \sigma_{e_{12}} \\ \sigma_{e_{12}} & \sigma_{e_2}^2 \end{bmatrix}$  is the residual variance-covariance matrix of the two traits.

167 Since STITCH provides for each SNP and each individual the imputed genotype  
 168 (the most likely genotype) as well as the genotype dosages (expected genotypic count,  
 169 weighted mean of the three genotype values, i.e., 2, 1, and 0, weighted by their  
 170 corresponding posterior probabilities), the  $\mathbf{G}$  matrix can be constructed using either the  
 171 imputed genotypes or the genotype dosages. The genotype-based  $\mathbf{G}$  matrix was  
 172 constructed with the method of VanRaden [22] as follows.

$$173 \mathbf{G} = \mathbf{W}\mathbf{W}' / \sum 2p_j(1 - p_j)$$

174 Here,  $\mathbf{W}$  is the centralized maker genotype matrix with its  $ij$ th element equal to

$$175 w_{ij} = m_{ij} - 2p_j$$

176 where  $m_{ij}$  (= 2,1, or 0) is the original genotype of individual  $i$  for SNP  $j$ ,  $p_j$  is the  
177 allele frequency of SNP  $j$ .

178 The dosage-based  $\mathbf{G}$  matrix was constructed as follows.

$$179 \quad \mathbf{G} = \mathbf{D}\mathbf{D}'/s_d$$

180 Here,  $\mathbf{D}$  is the centralized marker dosage matrix whose elements are zero-centered  
181 genotype dosages.  $s_d$  is the sum of variances for every column of  $\mathbf{D}$ .

182 We used both of the two strategies to construct  $\mathbf{G}$ . Different marker densities were  
183 also considered when constructing  $\mathbf{G}$  to evaluate the effect of marker density on the  
184 performance of genomic prediction.

185 We used GMAT [23] to construct the  $\mathbf{G}$  matrix and DMU (<http://dmu.agrsci.dk>) to  
186 estimate the variance and covariance components involved in the models and the  
187 genomic breeding values (GEBVs).

188

### 189 *Cross-validation*

190 In this study, a 12-fold Cross-validation (CV) was applied to assess the accuracy of  
191 the genomic selection. The 594 animals were divided into 12 subsets. One of them was  
192 taken in turn to be used as validation population and the rest 11 subsets used as training  
193 population. The accuracy of genomic prediction for the validation animals was  
194 evaluated by  $\frac{r_{y_c, GEBV}}{h}$ , i.e., the correlation between corrected phenotypic values ( $y_c$ ) and  
195 GEBVs divided by the square root of the heritability. The unbiasedness of predictions  
196 was assessed by the regression of  $y_c$  on GEBV ( $b_{y_c, GEBV}$ ).

197 The corrected phenotype for each animal was calculated as the original phenotypic

198 value corrected for fixed effects (sex and year-season effects) which were estimated  
199 using the conventional BLUP model based on the full dataset, i.e.,  $y_c = y - \text{sex effect} -$   
200 year-season effect.

201

## 202 **Results**

### 203 **Accuracies of genotype imputation**

#### 204 *Comparison of different pipelines*

205 The two genotype imputation pipelines, BaseVar + STITCH and Bcftools + Beagle,  
206 were compared using the original sequencing data of the 617 animals with average  
207 sequencing depth of 3.5x. It was obvious that the BaseVar + STITCH pipeline was  
208 significantly better than the Bcftools + Beagle pipeline (Figure 1), the average genotype  
209 accuracy from BaseVar + STITCH was about 6.5 percentage points higher than that  
210 from Bcftools + Beagle and the average genotype concordance was about 2.4  
211 percentage points higher. Therefore, the BaseVar + STITCH pipeline was adapted for  
212 the subsequent analyses.

213

#### 214 *The effects of sample size and sequencing depth*

215 We compared the genotype accuracy and genotype concordance for imputation  
216 with different sample sizes (200, 400 and 600) and sequencing depths (1x, 1.5x and  
217 3.5x) (Figure 2). In general, as expected, the genotype accuracy and genotype  
218 concordance increased with the increase of sample size and sequencing depth. The  
219 improvement of imputation accuracy was more obvious when the sample size was

220 increased from 200 to 400 and the sequencing depth increased from 1x to 1.5x. It should  
221 be noted that with sample size of  $\geq 400$  very high imputation accuracy (genotype  
222 accuracy  $> 0.94$ , genotype concordance  $> 0.98$ ) could be achieved even when the  
223 sequencing depth was as low as 1x.

224

### 225 *The effect of MAF*

226 Figure 3 shows the effect of MAF on imputation accuracy for sample size of 600.  
227 For SNPs with  $MAF < 0.01$ , the imputation accuracy was greatly affected by MAF and  
228 the accuracy increased rapidly with the increase of MAF. However, for SNPs with  
229  $MAF > 0.01$ , the imputation accuracy was not affected by MAF.

230

### 231 **Variance component estimation**

232 For the single-trait model, the variance components and heritabilities of BW and  
233 WW were estimated with the **G** matrix constructed using four levels of marker densities  
234 and two strategies (genotype-based **G** matrix and dosage-based **G** matrix). From the  
235 original sequence data with an average depth of 3.5x, we obtained 2.3M SNPs after  
236 imputation and quality control. We then reduced the marker density by applying LD  
237 pruning with three coefficients ( $r^2$ : 0.2, 0.4, and 0.8) by PLINK [19], leaving 130K,  
238 220K and 410K SNPs, respectively. The estimates under the four marker densities were  
239 very similar (Table 1). The estimates of additive genetic variances and heritabilities  
240 based on the dosage-based **G** matrix were all smaller than that based on the genotype-  
241 based **G** matrix, although the differences were very small.

242 For the two-trait model, the variance and co-variance components of the two traits  
243 were estimated based on the dosage-based **G** matrix constructed using 410K SNPs  
244 (Table 2). The estimates of heritability from the two-trait model (0.627 for birth weight,  
245 0.425 for weaning weight) were higher than that from single-trait model (0.580 for birth  
246 weight, 0.330 for weaning weight). The estimate of genetic correlation between birth  
247 weight and weaning weight was 0.839.

248

### 249 **Accuracy and unbiasedness of GEBVs**

250 The GEBVs for BW and WW were calculated under single-trait model and two-  
251 trait model, respectively. For single-trait model, we again considered different **G**  
252 matrices constructed with the four levels of marker densities and the two strategies as  
253 mentioned above. The accuracies and unbiasedness derived from cross-validation are  
254 given in Table 3. Marker density of 410K resulted slightly better results than the other  
255 marker densities. The results from the genotype-based **G** matrix and the dosage-based  
256 **G** matrix were very similar, with the accuracies from the latter very slightly higher than  
257 that from the former. For the two-trait model, only the dosage-based **G** matrix from  
258 410k markers was used. Compared with the results under single-trait model with the  
259 same **G** matrix, the two-trait model improved the accuracies by 4 and 5.5 percentage  
260 points for BW and WW, respectively, while maintained the same unbiasedness (Table  
261 4).

262

### 263 **Discussion**

264 Imputation-based lcWGS data is increasingly being used for genetic analysis of  
265 complex traits, such as genome-wide association study and genomic  
266 selection/prediction, in human [24-26], plants [27-28] and animals [10, 29-30], and has  
267 been proved to be a cost-effective way for genome-wide high-density genotyping,  
268 especially for species (such as donkey) for which a SNP array is not available.

269 Imputation is necessary for lcWGS data due to the high missing rates. Most of the  
270 proposed imputation methods designed for lcWGS data infer the gaps between the  
271 sparsely mapped reads by leveraging information from a reference panel of haplotypes  
272 [31-32]. However, for some animal species there is no such a reference panel available  
273 or the size of the panel is not large enough to provide sufficient and reliable information.  
274 Some imputation software, like Beagle, can work with or without a haplotype reference  
275 panel, however, the imputation accuracy with a haplotype reference panel is much  
276 better than that without a haplotype reference panel [12, 33]. Davies et al. (2016)  
277 developed a method, STITCH, which was designed specifically for the situation of  
278 without a haplotype reference panel, and was proved to yield accuracy comparable with  
279 Beagle with a haplotype reference panel [12].

280 An important step before imputation is to call SNPs from sequence data. GATK  
281 [34] and Bcftools<sup>[18]</sup> are two commonly used software for SNP calling from sequence  
282 data. Recently, Liu et al. (2018) develop a SNP calling method, BaseVar [7], which was  
283 designed for ultra-low coverage sequencing data. In this study, we compared two  
284 pipelines for SNP calling and imputation, i.e., Bcftools + Beagle and BaseVar +  
285 STITCH. We demonstrated that BaseVar + STITCH overperformed Bcftools + Beagle

286 (Figure 1). We showed that in our Dezhou Donkey population, using this pipeline, high  
287 imputation accuracy (genotype accuracy  $> 0.94$ , genotype concordance  $> 98\%$ ) can be  
288 achieved with a sample size of 400 and sequencing depth of 1x (Figure 2). Similar  
289 results were also reported by Zhang et al. [10]. In other words, with a sample size of  
290 over 400, a sequencing depth of 1x could be sufficient to ensure high imputation  
291 accuracy using BaseVar + STITCH.

292 Using the imputation-based sequence data, we evaluated the performance of  
293 genomic selection using GBLUP with respect to four different marker densities (130K,  
294 220K, 410K, and 2.3M), two different **G** matrix construction strategies (genotype-  
295 based **G** vs dosage-based **G**), and single-trait vs two-trait models. We found that the  
296 prediction accuracy increased slightly when the marker density increased from 130K to  
297 410K. However, it did not further increase when the density increased to 2.3M. The  
298 densities of 130K, 220K, and 410K correspond to medium to high density of SNP array,  
299 while the 2.3M density corresponds to the density of sequence data. It has been reported  
300 that, in the frame of GBLUP, the genomic prediction accuracy could be improved using  
301 high density SNP array compared to using medium density array [35-37], however,  
302 sequence data could hardly improve the accuracy compared with SNP array [36, 38].  
303 However, sequence data can be meaningful for cross-breed/population genomic  
304 selection [39-40]. The accuracy of genomic prediction using the dosage-based **G** matrix  
305 is only slightly better than that using the genotype-based **G** matrix. Since the  
306 improvement was rather small, we infer this may be due to the high accuracy of  
307 imputation.

308 Noticeable increases in genomic prediction accuracy were observed when using a  
309 two-trait model compared with using a single-trait model. It has been long widely  
310 proved that multi-trait model can increase the accuracy of breeding value estimation,  
311 either by conventional BLUP or by GBLUP [41-43], in particular for traits with high  
312 genetic correlation. For the two traits analyzed in this study, birth weight and weaning  
313 weight, we obtained an estimated genetic correlation of 0.839, which is very high. This  
314 increasement in accuracy with multi-trait model will be particularly beneficial for the  
315 situation where the reference population size is limited.

316

## 317 **Conclusions**

318 In this study we demonstrated that the pipeline BaseVar + STITCH is a good choice for  
319 SNP calling and imputation for low coverage sequence data. Sufficient high imputation  
320 accuracy could be achieved for sequence data with sequencing depth as low as 1x, when  
321 the size of the sequencing population is over 400. Thus, lcWGS combined with  
322 imputation provides a cost-effective way for whole genome high density genotyping  
323 and can be applied for large scale genomic selection in farm animals. This is particularly  
324 beneficial for those animal species for which a SNP array is not available. In frame of  
325 GBLUP, increasing marker density from a density comparable with a high-density SNP  
326 array (e.g., 400K) to sequence density with millions of SNPs did not increase the  
327 accuracy of genomic selection. Multi-trait model GBLUP improves significantly the  
328 accuracy of genomic selection over single-trait model, which would be particularly  
329 meaningful for the situation where the reference population size is limited.



## 330 **Acknowledgements**

331 The authors thank Supercomputing Center in Shandong Agricultural University for  
332 technical support.

## 333 **Authors' contributions**

334 QZ and CN designed the study. CZ, XZ, XJ, HL and SL collected the sample and  
335 performed the experiments. CZ and JT analyzed and interpreted the data. CZ, JT, XZ,  
336 CN, DW and QZ drafted the manuscript.

## 337 **Abbreviations**

338 GS: Genomic selection; WGS: Whole genome sequencing; LcWGS: Low coverage  
339 whole genome sequencing; MAF: minor allele frequency; BW: Birth weight; WW:  
340 weaning weight; GBLUP: Genomic best linear unbiased prediction; GEBV: Genomic  
341 breeding value; CV: Cross-validation

## 342 **Funding**

343 Project for Improved Agricultural Breeding of Shandong Province (2019LZGC011)

## 344 **Availability of data and materials**

345 All data supporting our findings are included in the manuscript.

## 346 **Ethics approval and consent to participate**

347 All experimental chickens were maintained, and all the studies were carried out  
348 according to the guideline of the experimental animal management of Shandong  
349 Agricultural University (SDAUA-2018-018).

## 350 **Consent for publication**

351 Not applicable.

352 **Competing interests**

353 The authors declare that they have no competing interests.

354 **References**

- 355 1. Meuwissen TH, Hayes BJ, Goddard ME. Prediction of total genetic value using genome-wide dense  
356 marker maps. *Genetics*. 2001;157(4): 1819-29.
- 357 2. Schaeffer LR. Strategy for applying genome-wide selection in dairy cattle. *J Anim Breed Genet*.  
358 2006;123(4): 218-23. doi:10.1111/j.1439-0388.2006.00595.x.
- 359 3. Wiggans GR, Cole JB, Hubbard SM, Sonstegard TS. Genomic selection in dairy cattle: The USDA  
360 experience. *Annu Rev Anim Biosci*. 2017;5(309-27). doi:10.1146/annurev-animal-021815-111422.
- 361 4. Yang AQ, Chen B, Ran ML, Yang GM, Zeng C. The application of genomic selection in pig cross  
362 breeding. *Yi Chuan*. 2020;42(2): 145-52. doi:10.16288/j.ycz.19-253.
- 363 5. van Binsbergen R, Calus MP, Bink MC, van Eeuwijk FA, Schrooten C, Veerkamp RF. Genomic  
364 prediction using imputed whole-genome sequence data in Holstein Friesian cattle. *Genet Sel Evol*.  
365 2015;47(71). doi:10.1186/s12711-015-0149-x.
- 366 6. Moghaddar N, Khansefid M, van der Werf J, Bolormaa S, Duijvesteijn N, Clark SA et al. Genomic  
367 prediction based on selected variants from imputed whole-genome sequence data in Australian  
368 sheep populations. *Genet Sel Evol*. 2019;51(1): 72. doi:10.1186/s12711-019-0514-2.
- 369 7. Liu S, Huang S, Chen F, Zhao L, Yuan Y, Francis SS et al. Genomic analyses from non-invasive  
370 prenatal testing reveal genetic associations, patterns of viral infections, and chinese population  
371 history. *Cell*. 2018;175(2): 347-59. doi:10.1016/j.cell.2018.08.016.
- 372 8. Nicod J, Davies RW, Cai N, Hassett C, Goodstadt L, Cosgrove C et al. Genome-wide association  
373 of multiple complex traits in outbred mice by ultra-low-coverage sequencing. *Nat Genet*.

- 374 2016;48(8): 912-8. doi:10.1038/ng.3595.
- 375 9. Pasaniuc B, Rohland N, McLaren PJ, Garimella K, Zaitlen N, Li H et al. Extremely low-coverage  
376 sequencing and imputation increases power for genome-wide association studies. *Nat Genet.*  
377 2012;44(6): 631-5. doi:10.1038/ng.2283.
- 378 10. Zhang W, Li W, Liu G, Gu L, Ye K, Zhang Y et al. Evaluation for the effect of low-coverage  
379 sequencing on genomic selection in large yellow croaker. *Aquaculture.* 2021;534(
- 380 11. Hickey JM. Sequencing millions of animals for genomic selection 2.0. *J Anim Breed Genet.*  
381 2013;130(5): 331-2. doi:10.1111/jbg.12054.
- 382 12. Davies RW, Flint J, Myers S, Mott R. Rapid genotype imputation from sequence without reference  
383 panels. *Nat Genet.* 2016;48(8): 965-9. doi:10.1038/ng.3594.
- 384 13. Hui R, D'Atanasio E, Cassidy LM, Scheib CL, Kivisild T. Evaluating genotype imputation pipeline  
385 for ultra-low coverage ancient genomes. *Sci Rep.* 2020;10(1): 18542. doi:10.1038/s41598-020-  
386 75387-w.
- 387 14. Ros-Freixedes R, Gonen S, Gorjanc G, Hickey JM. A method for allocating low-coverage  
388 sequencing resources by targeting haplotypes rather than individuals. *Genet Sel Evol.* 2017;49(1):  
389 78. doi:10.1186/s12711-017-0353-y.
- 390 15. Wang C, Li H, Guo Y, Huang J, Sun Y, Min J et al. Donkey genomes provide new insights into  
391 domestication and selection for coat color. *Nat Commun.* 2020;11(1): 6014. doi:10.1038/s41467-  
392 020-19813-7.
- 393 16. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform.  
394 *Bioinformatics.* 2009;25(14): 1754-60. doi:10.1093/bioinformatics/btp324.
- 395 17. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N et al. The Sequence Alignment/Map

- 396 format and SAMtools. *Bioinformatics*. 2009;25(16): 2078-9. doi:10.1093/bioinformatics/btp352.
- 397 18. Danecek P, Bonfield JK, Liddle J, Marshall J, Ohan V, Pollard MO et al. Twelve years of SAMtools  
398 and BCFtools. *Gigascience*. 2021;10(2). doi:10.1093/gigascience/giab008.
- 399 19. Chang CC, Chow CC, Tellier LC, Vattikuti S, Purcell SM, Lee JJ. Second-generation PLINK:  
400 Rising to the challenge of larger and richer datasets. *Gigascience*. 2015;4(7). doi:10.1186/s13742-  
401 015-0047-8.
- 402 20. Browning BL, Browning SR. Genotype imputation with millions of reference samples. *Am J Hum*  
403 *Genet*. 2016;98(1): 116-26. doi:10.1016/j.ajhg.2015.11.020.
- 404 21. Habier D, Fernando RL, Dekkers JC. The impact of genetic relationship information on genome-  
405 assisted breeding values. *Genetics*. 2007;177(4): 2389-97. doi:10.1534/genetics.107.081190.
- 406 22. VanRaden PM. Efficient methods to compute genomic predictions. *J Dairy Sci*. 2008;91(11): 4414-  
407 23. doi:10.3168/jds.2007-0980.
- 408 23. Wang D, Tang H, Liu JF, Xu S, Zhang Q, Ning C. Rapid epistatic mixed-model association studies  
409 by controlling multiple polygenic effects. *Bioinformatics*. 2020;36(19): 4833-7.  
410 doi:10.1093/bioinformatics/btaa610.
- 411 24. Rustagi N, Zhou A, Watkins WS, Gedvilaite E, Wang S, Ramesh N et al. Extremely low-coverage  
412 whole genome sequencing in South Asians captures population genomics information. *BMC*  
413 *Genomics*. 2017;18(1): 396. doi:10.1186/s12864-017-3767-6.
- 414 25. Homburger JR, Neben CL, Mishne G, Zhou AY, Kathiresan S, Khera AV. Low coverage whole  
415 genome sequencing enables accurate assessment of common variants and calculation of genome-  
416 wide polygenic scores. *Genome Med*. 2019;11(1): 74. doi:10.1186/s13073-019-0682-2.
- 417 26. Pasaniuc B, Rohland N, McLaren PJ, Garimella K, Zaitlen N, Li H et al. Extremely low-coverage

- 418 sequencing and imputation increases power for genome-wide association studies. *Nat Genet.*  
419 2012;44(6): 631-5. doi:10.1038/ng.2283.
- 420 27. Gardner EM, Johnson MG, Ragone D, Wickett NJ, Zerega NJ. Low-coverage, whole-genome  
421 sequencing of *Artocarpus camansi* (Moraceae) for phylogenetic marker development and gene  
422 discovery. *Appl Plant Sci.* 2016;4(7). doi:10.3732/apps.1600017.
- 423 28. Zhou C, Duarte T, Silvestre R, Rossel G, Mwanga R, Khan A et al. Insights into population structure  
424 of East African sweetpotato cultivars from hybrid assembly of chloroplast genomes. *Gates Open*  
425 *Res.* 2018;2(41). doi:10.12688/gatesopenres.12856.2.
- 426 29. Yang R, Guo X, Zhu D, Bian C, Zhao Y, Tan C et al. Genome-wide association analyses of multiple  
427 traits in Duroc pigs using low-coverage. *bioRxiv.* 2019. doi:10.1101/754671.
- 428 30. Xu C, Wu K, Zhang JG, Shen H, Deng HW. Low-, high-coverage, and two-stage DNA sequencing  
429 in the design of the genetic association study. *Genet Epidemiol.* 2017;41(3): 187-97.  
430 doi:10.1002/gepi.22015.
- 431 31. Rubinacci S, Ribeiro DM, Hofmeister RJ, Delaneau O. Efficient phasing and imputation of low-  
432 coverage sequencing data using large reference panels. *Nat Genet.* 2021;53(1): 120-6.  
433 doi:10.1038/s41588-020-00756-0.
- 434 32. Spiliopoulou A, Colombo M, Orchard P, Agakov F, McKeigue P. GeneImp: Fast imputation to  
435 large reference panels using genotype likelihoods from ultralow coverage sequencing. *Genetics.*  
436 2017;206(1): 91-104. doi:10.1534/genetics.117.200063.
- 437 33. Zheng C, Boer MP, van Eeuwijk FA. Accurate Genotype Imputation in Multiparental Populations  
438 from Low-Coverage Sequence. *Genetics.* 2018;210(1): 71-82. doi:10.1534/genetics.118.300885.
- 439 34. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytzky A et al. The Genome

- 440 Analysis Toolkit: A MapReduce framework for analyzing next-generation DNA sequencing data.  
441 Genome Res. 2010;20(9): 1297-303. doi:10.1101/gr.107524.110.
- 442 35. Boison SA, Utsunomiya A, Santos D, Neves H, Carvalheiro R, Meszaros G et al. Accuracy of  
443 genomic predictions in Gyr (*Bos indicus*) dairy cattle. J Dairy Sci. 2017;100(7): 5479-90.  
444 doi:10.3168/jds.2016-11811.
- 445 36. Perez-Enciso M, Rincon JC, Legarra A. Sequence- vs. Chip-assisted genomic selection: Accurate  
446 biological information is advised. Genet Sel Evol. 2015;47(43). doi:10.1186/s12711-015-0117-5.
- 447 37. VanRaden PM, Null DJ, Sargolzaei M, Wiggans GR, Tooker ME, Cole JB et al. Genomic  
448 imputation and evaluation using high-density Holstein genotypes. J Dairy Sci. 2013;96(1): 668-78.  
449 doi:10.3168/jds.2012-5702.
- 450 38. Perez-Enciso M. Genomic relationships computed from either next-generation sequence or array  
451 SNP data. J Anim Breed Genet. 2014;131(2): 85-96. doi:10.1111/jbg.12074.
- 452 39. Druet T, Macleod IM, Hayes BJ. Toward genomic prediction from whole-genome sequence data:  
453 Impact of sequencing design on genotype imputation and accuracy of predictions. Heredity (Edinb).  
454 2014;112(1): 39-47. doi:10.1038/hdy.2013.13.
- 455 40. van Binsbergen R, Bink MC, Calus MP, van Eeuwijk FA, Hayes BJ, Hulsegge I et al. Accuracy of  
456 imputation to whole-genome sequence data in Holstein Friesian cattle. Genet Sel Evol. 2014;46(41).  
457 doi:10.1186/1297-9686-46-41.
- 458 41. Guo G, Zhao F, Wang Y, Zhang Y, Du L, Su G. Comparison of single-trait and multiple-trait  
459 genomic prediction models. BMC Genet. 2014;15(30). doi:10.1186/1471-2156-15-30.
- 460 42. Jia Y, Jannink JL. Multiple-trait genomic selection methods increase genetic value prediction  
461 accuracy. Genetics. 2012;192(4): 1513-22. doi:10.1534/genetics.112.144246.

462 43. Calus MP, Veerkamp RF. Accuracy of multi-trait genomic selection using different methods. *Genet*  
463 *Sel Evol.* 2011;43(26). doi:10.1186/1297-9686-43-26.

464

## 465 **Figures**

### 466 **Figure 1 Comparison of imputation accuracy using two imputation pipelines**

467 BaseVar + STITCH: Calling SNPs using BaseVar and imputing missing genotypes using STITCH

468 Bcftools + Beagle: Calling SNPs using Bcftools and imputing missing genotypes using Beagle

469 a, b and c: genotype accuracy for chromosomes 1, 19, and 30, respectively;

470 d, e and f: genotype concordance for chromosomes 1, 19 and 30, respectively

471

### 472 **Figure 2 The effects of sample size and sequencing depth on imputation accuracy** 473 **using the pipeline BaseVar + STITCH**

474 a, b and c: genotype accuracy for chromosomes 1, 19, and 30, respectively;

475 d, e and f: genotype concordance for chromosomes 1, 19, and 30, respectively

476

### 477 **Figure 3 The effect of minor allele frequency on imputation accuracy using the** 478 **pipeline BaseVar + STITCH**

479 a, b and c: genotype accuracy for chromosomes 1, 19, and 30, respectively;

480 d, e and f: genotype concordance for chromosomes 1, 19, and 30, respectively

481

482

483

484 **Table 1** Estimates of variance components and heritabilities based on single-trait model  
 485 for birth weight (BW) and weaning weight (WW) in Dezhou Donkey population

Variance	130K		220K		410K		2.3M	
	G(g)	G(d)	G(g)	G(d)	G(g)	G(d)	G(g)	G(d)
BW								
$\sigma_a^2$	10.266	10.100	10.424	10.253	10.750	10.581	10.166	10.011
$\sigma_e^2$	8.024	8.030	7.904	7.912	7.656	7.664	8.154	8.159
$h^2$	0.561	0.557	0.569	0.564	0.584	0.580	0.555	0.551
WW								
$\sigma_a^2$	68.419	68.223	69.136	68.919	69.866	69.670	62.042	61.083
$\sigma_e^2$	139.977	140.874	140.659	141.777	141.046	141.328	146.029	146.074
$h^2$	0.328	0.326	0.330	0.327	0.331	0.330	0.298	0.295

486  $\sigma_a^2$ : additive genetic variance;  $\sigma_e^2$ : residual variance;  $h^2$ : heritability; G(g): genotype-based G  
 487 matrix; G(d): dosage-based G matrix

488

489 **Table 2** Estimates of variance components, heritabilities, and genetic correlation for  
 490 birth weight (BW) and weaning weight (WW) under a two-trait model

Variance	Two-trait model <sup>a</sup>	
	BW	WW
$\sigma_a^2$	11.769	90.728
$\sigma_e^2$	7.016	122.553
$h^2$	0.627	0.425
$r_g$	0.839	

491 a: using dosage-based G matrix constructed with 410K SNPs

492  $\sigma_a^2$ : additive genetic variance;  $\sigma_e^2$ : residual variance;  $h^2$ : heritability;  $r_g$ : genetic correlation;

493



494 **Table 3** Accuracies ( $\pm SEs$ ) and unbiasedness ( $\pm SEs$ ) of genomic selection under a  
 495 single-trait model for birth weight (BW) and weaning weight (WW)

Marker density	Genotype-based <b>G</b> matrix		Dosage-based <b>G</b> matrix	
	Accuracy	Unbiasedness	Accuracy	Unbiasedness
BW				
130K	0.368 $\pm$ 0.042	0.983 $\pm$ 0.141	0.370 $\pm$ 0.041	0.983 $\pm$ 0.141
220K	0.370 $\pm$ 0.042	0.984 $\pm$ 0.139	0.371 $\pm$ 0.040	0.984 $\pm$ 0.140
410K	0.376 $\pm$ 0.041	0.990 $\pm$ 0.135	0.378 $\pm$ 0.040	0.990 $\pm$ 0.136
2.3M	0.371 $\pm$ 0.044	0.980 $\pm$ 0.140	0.372 $\pm$ 0.042	0.981 $\pm$ 0.140
WW				
130K	0.399 $\pm$ 0.030	1.198 $\pm$ 0.209	0.400 $\pm$ 0.030	1.199 $\pm$ 0.210
220K	0.402 $\pm$ 0.030	1.191 $\pm$ 0.213	0.403 $\pm$ 0.029	1.191 $\pm$ 0.209
410K	0.409 $\pm$ 0.029	1.187 $\pm$ 0.214	0.410 $\pm$ 0.029	1.187 $\pm$ 0.214
2.3M	0.403 $\pm$ 0.032	1.212 $\pm$ 0.221	0.405 $\pm$ 0.031	1.212 $\pm$ 0.221

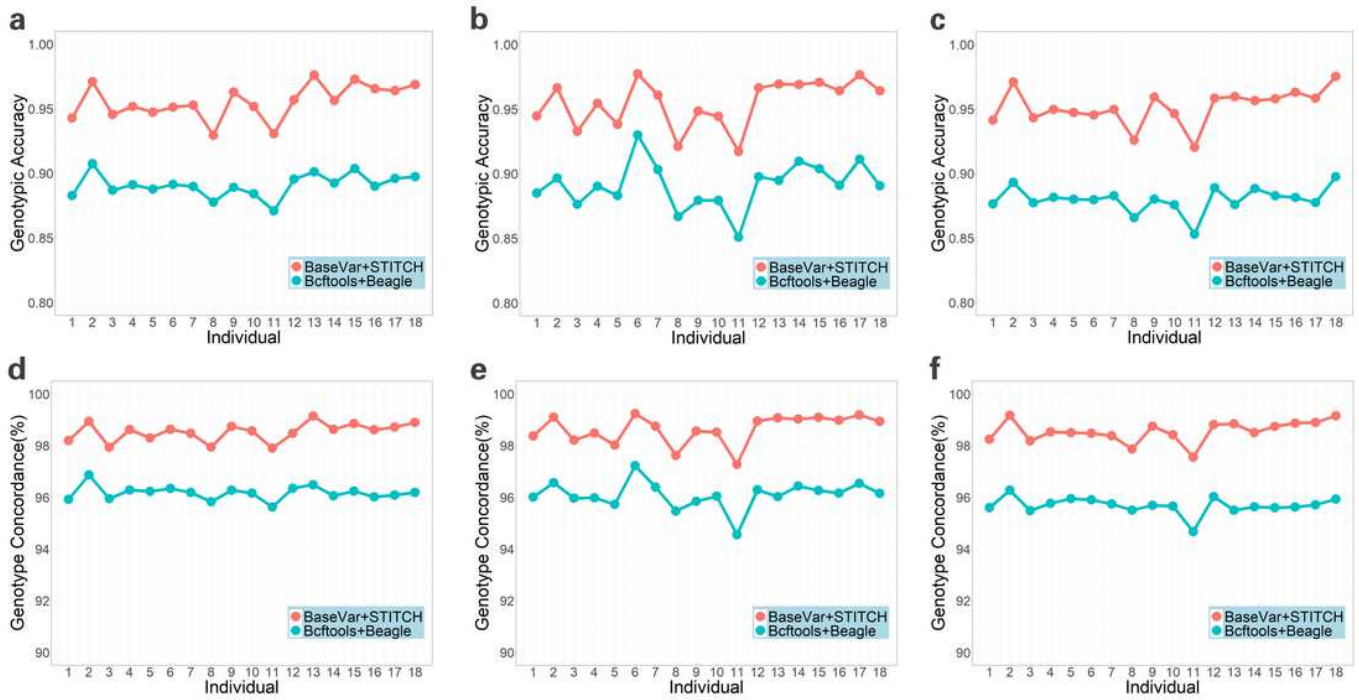
496

497 **Table 4** Accuracies ( $\pm SEs$ ) and unbiasedness ( $\pm SEs$ ) of genomic selection for birth  
 498 weight (BW) and weaning weight (WW) under a single- and a two-trait model using  
 499 dosage-based **G** matrix constructed with 410K SNPs

Model	BW		WW	
	Accuracy	Unbiasedness	Accuracy	Unbiasedness
Single-trait model	0.378 $\pm$ 0.041	0.990 $\pm$ 0.136	0.410 $\pm$ 0.029	1.187 $\pm$ 0.214
Two-trait model	0.417 $\pm$ 0.039	0.991 $\pm$ 0.113	0.465 $\pm$ 0.039	1.182 $\pm$ 0.153

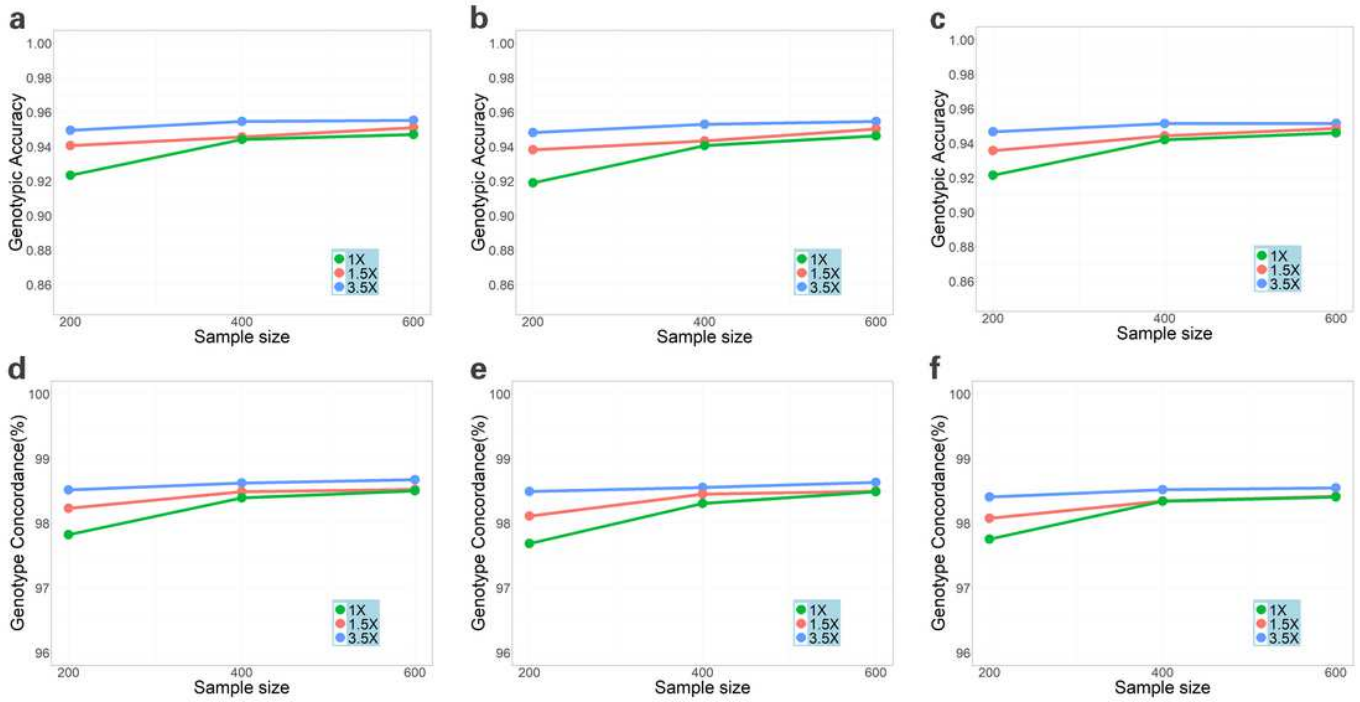
500

# Figures



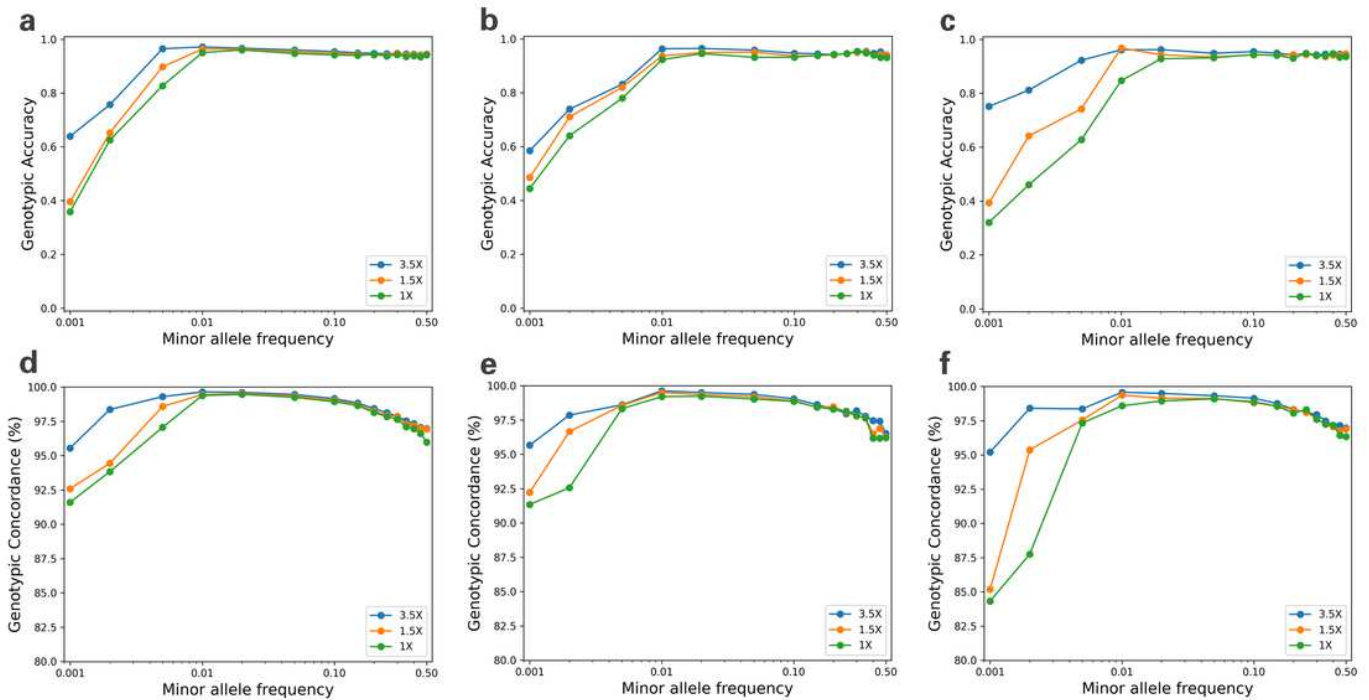
**Figure 1**

Comparison of imputation accuracy using two imputation pipelines BaseVar + STITCH: Calling SNPs using BaseVar and imputing missing genotypes using STITCH Bcftools + Beagle: Calling SNPs using Bcftools and imputing missing genotypes using Beagle a, b and c: genotypic accuracy for chromosomes 1, 19, and 30, respectively; d, e and f: genotype concordance for chromosomes 1, 19 and 30, respectively



**Figure 2**

The effects of sample size and sequencing depth on imputation accuracy using the pipeline BaseVar + STITCH a, b and c: genotype accuracy for chromosomes 1, 19, and 30, respectively; d, e and f: genotype concordance for chromosomes 1, 19, and 30, respectively



**Figure 3**

The effect of minor allele frequency on imputation accuracy using the pipeline BaseVar + STITCH a, b and c: genotype accuracy for chromosomes 1, 19, and 30, respectively; d, e and f: genotype concordance for chromosomes 1, 19, and 30, respectively