1     **JAST (Journal of Animal Science and Technology) TITLE PAGE**

2     **Upload this completed form to website with submission**

3

| ARTICLE INFORMATION | Fill in information in each box below |
|---|---|
| **Article Type** | Research Article |
| **Article Title (within 20 words without abbreviations)** | Effect of Breed Composition in Genomic Prediction Using Crossbred Pig Reference Population |
| **Running Title (within 10 words)** | Effect of Breed Composition in Genomic Prediction |
| **Author** | Euiseo Hong [first_author] 1, Yoonji Chung [first_author] 2, Phuong Thanh N. Dinh3, Yoonsik Kim2, Suyeon Maeng4, Young jae Choi3, Jaeho Lee3, Woonyoung Jeong1, Hyunji Choi5, Seung Hwan Lee4 |
| **Affiliation** | 1 Department of Bio-Big Data and Precision Agriculture, Chungnam National University, Daejeon 34134, Korea, Republic of<br>2 Institute of Agricultural Science, Chungnam National University, Daejeon 34134, Korea, Republic of<br>3 Department of Bio-AI Convergence, Chungnam National University, Daejeon 34134, Korea, Republic of<br>4 Division of Animal & Dairy Science, Chungnam National University, Daejeon 34134, Korea, Republic of<br>5 Division of Animal Genomics and Bioinformatics, National Institute of Animal Science, Jeonbuk-do 55365, Korea, Republic of |
| **ORCID (for more information, please visit https://orcid.org)** | Euiseo Hong (https://orcid.org/0000-0003-3078-2560)<br>Yoonji Chung (https://orcid.org/0000-0002-6906-6468)<br>Phuong Thanh N. Dinh (https://orcid.org/0000-0002-3057-0210)<br>Yoonsik Kim (https://orcid.org/0000-0002-5318-7521)<br>Suyeon Maeng (https://orcid.org/0000-0001-9903-3803)<br>Young jae Choi (https://orcid.org/0000-0003-1540-6970)<br>Jaeho Lee (https://orcid.org/0009-0008-7721-8135)<br>Woonyoung Jeong (https://orcid.org/0009-0002-7572-1382)<br>Hyunji Choi (https://orcid.org/0000-0001-9782-6586)<br>Seung Hwan Lee (https://orcid.org/0000-0003-1508-4887) |
| **Competing interests** | No potential conflict of interest relevant to this article was reported. |
| **Funding sources**<br>State funding sources (grants, funding sources, equipment, and supplies). Include name and number of grant if available. | *This work was supported by Chungnam National University.* |
| **Acknowledgements** | |
| **Availability of data and material** | |
| **Authors' contributions**<br>Please specify the authors' role using this form. | Conceptualization: Hong E, Lee SH<br>Data curation: Hong E, Chung Y, Choi H, Lee SH<br>Formal analysis: Hong E, Chung Y, Lee SH<br>Methodology: Hong E, Chung Y, Dinh PTN, Kim Y<br>Software: Hong E, Maeng S, Choi YJ, Lee J, Jeong W<br>Validation: Hong E<br>Investigation: Hong E<br>Writing - original draft: Hong E<br>Writing - review & editing: Hong E, Chung Y, Dinh PTN, Kim Y, Maeng S, Choi YJ, Lee J, Jeong W, Choi H, Lee SH |

| Ethics approval and consent to participate | This article does not require IRB/IACUC approval because there are no human and animal participants. |
|---|---|

4

## 5 CORRESPONDING AUTHOR CONTACT INFORMATION

| For the corresponding author (responsible for correspondence, proofreading, and reprints) | Fill in information in each box below |
|---|---|
| First name, middle initial, last name | Seung Hwan Lee |
| Email address – this is where your proofs will be sent | genomicselection46@gmail.com |
| Secondary Email address | |
| Address | |
| Cell phone number | +82-10-5337-1617 |
| Office phone number | +82-42-821-5878 |
| Fax number | |

6

7

8

# Effect of Breed Composition in Genomic Prediction Using Crossbred Pig Reference Population

## Abstract

In contrast to conventional genomic prediction, which typically targets a single breed and circumvents the necessity for population structure adjustments, multi-breed genomic prediction necessitates accounting for population structure to mitigate potential bias. The presence of this structure in multi-breed datasets can influence prediction accuracy, rendering proper modeling crucial for achieving unbiased results. This study aimed to address the effect of population structure on multi-breed genomic prediction, particularly focusing on crossbred reference populations. The predictive accuracy of genomic models was assessed by incorporating genomic breed composition (GBC) or principal component analysis (PCA) into the genomic best linear unbiased prediction (GBLUP) model. The accuracy of five different genomic prediction models was evaluated using data from 354 Duroc × Korean native pig crossbreds, 1,105 Landrace × Korean native pig crossbreds, and 1,107 Landrace × Yorkshire × Duroc crossbreds. The models tested were GBLUP without population structure adjustment, GBLUP with PCA as a fixed effect, GBLUP with GBC as a fixed effect, GBLUP with PCA as a random effect, and GBLUP with GBC as a random effect. The highest predictive accuracies for backfat thickness (0.59) and carcass weight (0.50) were observed in Models 1, 4, and 5. In contrast, Models 2 and 3, which included population structure as a fixed effect, exhibited lower accuracies, with backfat thickness accuracies of 0.40 and 0.53 and carcass weight accuracies of 0.34 and 0.38, respectively. These findings suggest that in multi-breed genomic prediction, the most efficient and accurate approach is either to forgo adjusting for population structure or, if adjustments are necessary, to model it as a random effect. This study provides a robust framework for multi-breed genomic prediction, highlighting the critical role of appropriately accounting for population structure. Moreover, our findings have important implications for improving genomic selection efficiency, ultimately enhancing commercial production by optimizing prediction accuracy in crossbred populations.

**Keywords: genomic breed composition, genomic prediction, multi-breed genomic prediction, population structure**

# Introduction

Accurate prediction of genomic breeding values is a critical component of successful genomic selection, which requires a sufficiently large reference population to reliably estimate marker effects [1]. However, small populations, such as Jersey cattle, often pose challenges owing to the limited reference populations of progeny-tested bulls, leading to less reliable genomic breeding values [2]. Consequently, genetic progress is restricted in breeds without a large reference population. One approach to addressing this limitation is across-breed prediction, which involves the use of a large reference dataset from another breed [3]. Another approach is multi-breed prediction, which combines data from multiple breeds to create a larger, more comprehensive dataset [3]. Both approaches can enhance prediction accuracy for smaller breeds, helping them become more competitive while minimizing the additional costs associated with genotyping and phenotyping.

Empirical studies have demonstrated that the accuracy of across-breed genomic prediction is often near zero and that combining multiple breeds has not yielded significant improvements in accuracy [3, 4]. However, these methods remain promising, particularly when combined with strategies that account for population structure and other sources of variation [5, 6]. Addressing population structure, also referred to as population stratification, is critical for genomic prediction across different breeds. Population structure arises from differences in allele frequencies between subpopulations, which may result from geographic separation, or natural or artificial selection [7]. These differences can lead to spurious marker-trait associations [8, 9], potentially inflating estimates of genomic heritability [10] and introducing bias into genomic prediction accuracy [6].

To mitigate the effects of population structure, it is important to model it appropriately within genomic prediction models, particularly when combining data from multiple breeds. A common method involves incorporating principal components (PCs) derived from genomic data as a fixed effect in the prediction model [7]. However, incorporating PCs as a fixed effect can result in over-correction, as these components are derived from the genomic relationship matrix used in genomic prediction [11]. To address this limitation, in this study, PCs were modeled as a random effect to capture population structure without confounding the genomic relationship matrix. The predictive accuracy of these models was compared with those of models in which PCs were excluded. Additionally, breed composition, another explanatory factor for population structure, was modeled as either a fixed or random effect to adjust for population structure.

In this study, we evaluated the accuracy of genomic predictions using models that incorporated breed

63  composition and PCs as fixed and random effects and compared the results with those of a baseline model. This

64  study aimed to determine whether accounting for population structure using breed composition or PCs can

65  improve genomic prediction accuracy. The findings of this study may provide valuable insights into optimizing

66  genomic prediction models for populations with complex or diverse structures.

67

# Materials and Methods

69  **Animals, genotypes, and phenotypes**

70  The genotype dataset comprised data from 354 Duroc × Korean native pigs (DK), 1,105 Landrace × Korean

71  native pigs (LK), 1,017 Landrace × Yorkshire × Duroc (LYD) crossbreds, along with purebred animals. Crossbred

72  individuals were genotyped using the Illumina PorcineSNP60 Genotyping BeadChip, whereas genotype data for

73  purebred animals were provided by the Centre for Research in Agricultural Genomics [12]. Genotype data for the

74  Korean native pigs (KNPs) among the purebreds were provided by the National Institute of Animal Science in

75  Korea. Details regarding the number of animals, single nucleotide polymorphisms (SNPs), and average observed

76  heterozygosity rate for each breed are presented in **Supplementary Table 1**. The quality control process involved

77  the exclusion of SNPs located on sex chromosomes, with a genotype call rate below 90%, and with a minor allele

78  frequency below 1%. After merging datasets and applying the quality control process, a common set of 24,118

79  SNPs were retained for analysis.

80  Phenotypic data revealed differences in backfat thickness and carcass weight among the breeds. The LYD breed

81  exhibited the lowest backfat thickness, whereas the DK breed had the highest backfat thickness. Conversely, the

82  DK breed exhibited the lowest carcass weight, whereas LYD had the highest carcass weight. The carcass

83  performance of the breeds crossed with the KNP was lower than that of LYD. This finding aligns with the known

84  characteristics of the KNP breed, which is known for its good meat quality but poor growth rate [13]. Statistical

85  details for the phenotypes are provided in **Supplementary Table 2**.

86

87  **Principal Component Analysis**

88  Principal component analysis (PCA) was employed to investigate genetic differences between populations and

89  to correct for population structure. PCA simplifies data complexity while maintaining the underlying relationships

90  among the data points. When applied to biallelic genotype data, PCA identifies the eigenvalues and eigenvectors

91  of the covariance matrix of allele frequencies, thereby reducing the data to a limited number of dimensions known

92  as PCs. Each PC represents a proportion of the total genomic variation. Subsequently, the data are mapped onto

93  the space defined by these PC axes, facilitating the visualization of samples and their distances from each other

94  in a scatter plot. In this visualization, sample overlap indicates shared genetic identity, reflecting common ancestry

95  or origin [14].

96

97  **Genomic Breed Composition**

98  Genomic breed composition (GBC) was estimated from genomic data using a maximum likelihood model

99  implemented in ADMIXTURE v1.3.0 [15]. ADMIXTURE uses genotype data to cluster individuals into

100  subgroups based on a predetermined number of groups. The projection extension of the ADMIXTURE program

101  allows for estimating ancestry using predefined ancestral population allele frequencies. This extension enables

102  efficient ancestry inference across large genomic datasets, leveraging allele frequencies from reference panels,

103  such as the 1000 Genomes Project. Additionally, the projection approach is particularly advantageous for datasets

104  with significant population distribution imbalances, as such imbalances can adversely affect the accuracy of

105  ancestry inference [16].

106  The projection extension of the ADMIXTURE program was used to analyze the dataset owing to the imbalance

107  between purebred and crossbred samples. Ancestral population allele frequencies were estimated using the

108  purebred samples, whereas the GBC values of the crossbreds were estimated using the allele frequencies of the

109  purebreds.

110

111  **Statistical Models**

112  First, PCs and GBCs were calculated for each individual, which were subsequently used in five models to

113  predict genomic estimated breeding values (GEBV). Although additional fixed effects such as age and farm were

114  considered, age information was unavailable, and farm data showed high multicollinearity with the PC and GBC

115  values, which precluded their inclusion.

116     ·     Model 1 (NULL) is defined as follows:

117     $$\mathbf{y} = \mathbf{Xb} + \mathbf{Zg} + \mathbf{e},$$

118     where $\mathbf{y}$ represents the vector of trait records (backfat thickness or carcass weight); $\mathbf{b}$ indicates the vector of

119     fixed effects, including sex; $\mathbf{X}$ denotes the design matrix linking fixed effects to the records; $\mathbf{g}$ represents the

120     vector of random genetic effects, modeled as $\sim N\left(0, \mathbf{G}\sigma_g^2\right)$, with $\mathbf{G}$ being the genomic relationship matrix and $\sigma_g^2$

121     being the genetic variance captured by the SNPs; $\mathbf{Z}$ indicates the design matrix linking records to animals; and $\mathbf{e}$

122     denotes the vector of random deviations, modeled as $\sim N(0, \mathbf{I}\sigma_e^2)$, with $\mathbf{I}$ as an animal-by-animal identity matrix

123     and $\sigma_e^2$ representing the error variance. The GEBV for this model was predicted as $\mathbf{GEBV} = \hat{\mathbf{g}}$. The genomic

124     relationship matrix was constructed using GCTA v1.94.1 software according to the following equation [17]:

125     $$\mathbf{G}_{jk} = \frac{1}{N} \sum_{i=1}^{N} \frac{(x_{ij} - 2p_i)(x_{ik} - 2p_i)}{2p_i(1 - p_i)},$$

126     where $x_{ij}$ and $x_{ik}$ represent the genotypes (coded as 0, 1, or 2) of individuals $j$ and $k$ at SNP $i$. $p_i$ indicates the

127     allele frequency of SNP $i$, and $N$ denotes the total number of SNPs. The distribution of the diagonal and off-

128     diagonal elements of the genomic relationship matrix is shown in **Supplementary Figure 1**. The mean of the

129     diagonal elements is 1.03, indicating low inbreeding within the population. The mean of the off-diagonal elements

130     is 0, showing that individuals are genetically independent of each other.

131     ·     Model 2 (PC_F) is defined as follows:

132     $$\mathbf{y} = \mathbf{Xb} + \mathbf{Zg} + \mathbf{e},$$

133     where $\mathbf{y}$ represents the vector of trait records; $\mathbf{b}$ denotes the vector of fixed effects, which includes PC values (20

134     PCs) and sex; $\mathbf{X}$ indicates the design matrix linking fixed effects to records; $\mathbf{g}$ represents the vector of random

135     genetic effects; $\mathbf{Z}$ denotes the design matrix linking records to animals; and $\mathbf{e}$ indicates the vector of random

136     deviations. For this model, $\mathbf{GEBV} = \hat{\mathbf{g}}$.

137     ·     Model 3 (GBC_F) is defined as follows:

138     $$\mathbf{y} = \mathbf{Xb} + \mathbf{Zg} + \mathbf{e},$$

139     where $\mathbf{y}$ represents the vector of trait records; $\mathbf{b}$ denotes the vector of fixed effects, which includes GBC values

140     and sex (here, breed composition values represent the proportion of each individual's genome derived from the

141    four breeds: Duroc, KNP, Landrace, and Yorkshire); $\mathbf{X}$ indicates the design matrix linking fixed effects to records;

142    $\mathbf{g}$ represents the vector of random genetic effects; $\mathbf{Z}$ denotes the design matrix linking records to animals; and $\mathbf{e}$

143    indicates the vector of random deviations. For this model, $\mathbf{GEBV} = \hat{\mathbf{g}}$.

144    ·    Model 4 (PC_R) is defined as follows:

145    $$\mathbf{y} = \mathbf{Xb} + \mathbf{Zg} + \mathbf{Zpc} + \mathbf{e},$$

146    where $\mathbf{y}$ indicates the vector of trait records; $\mathbf{b}$ represents the vector of fixed effects, including sex; $\mathbf{X}$ denotes the

147    design matrix linking fixed effects to records; $\mathbf{g}$ indicates the vector of random genetic effects; $\mathbf{pc}$ denotes the

148    vector of random variables representing groups of PC values, which were clustered using the Gaussian Mixture

149    Model implemented in the 'mclust' R package [18]; $\mathbf{Z}$ indicates the design matrix linking records to animals; and

150    $\mathbf{e}$ denotes the vector of random deviations. For this model, $\mathbf{GEBV} = \hat{\mathbf{g}} + \widehat{\mathbf{pc}}$.

151    ·    Model 5 (GBC_R) is defined as follows:

152    $$\mathbf{y} = \mathbf{Xf} + \mathbf{Zg} + \mathbf{Zgbc} + \mathbf{e},$$

153    where $\mathbf{y}$ represents the vector of trait records; $\mathbf{b}$ denotes the vector of fixed effects, including sex; $\mathbf{X}$ indicates

154    the design matrix linking fixed effects to records; $\mathbf{g}$ represents the vector of random genetic effects; $\mathbf{gbc}$ denotes

155    the vector of random variables representing groups of GBC values, which were clustered using the Gaussian

156    Mixture Model implemented in the 'mclust' R package [18]; $\mathbf{Z}$ indicates the design matrix linking records to

157    animals; and $\mathbf{e}$ represents the vector of random deviations. For this model, $\mathbf{GEBV} = \hat{\mathbf{g}} + \widehat{\mathbf{gbc}}$.

158    Variance components were estimated using the restricted maximum likelihood (REML) method, as

159    implemented in MTG2 [19], for each model. Heritability for the traits was estimated using the formula $h^2 =$

160    $\widehat{\sigma_g^2}/(\widehat{\sigma_g^2} + \widehat{\sigma_e^2})$. The accuracy of GEBVs for each of the five models was calculated as $r(\mathbf{GEBV}, \mathbf{y})$, where $\mathbf{y}$

161    represents the phenotypes corrected for fixed effects [20]. A 5-fold cross-validation approach was used to validate

162    the models. In this method, animals were randomly divided into five groups, with each group treated as the

163    validation set while the remaining groups constituted the reference set.

164

165    # Results

**Principal Components Analysis**

PCA was performed to explore genetic structure across populations. The analysis revealed that the first PC (PC1) accounted for 43.9% of the total genetic variance, whereas the second PC (PC2) constituted 13.6% of the variance (**Figure 1**). The PCA plot revealed a clear separation among the crossbred populations, indicating distinct genetic backgrounds. However, the LYD population exhibited greater dispersion along the first two PCs, suggesting more considerable genetic variation within this group. This observed variation is likely attributed to the presence of F1 hybrids in the dataset, which primarily combined Landrace and Yorkshire genetics, thereby increasing the overall diversity observed in this population.

**Genomic Breed Composition**

The breed composition of the crossbred populations was evaluated using ADMIXTURE analysis; the results are depicted in **Figure 2**. The analysis was conducted in unsupervised mode using genomic data from purebred samples, and the estimated breed allele frequencies were subsequently used to infer breed membership coefficients for the crossbred individuals.

In the LYD population, the estimated breed composition revealed an average contribution of 31%, 33%, and 36% from Landrace, Yorkshire, and Duroc, respectively (**Table 1**). The presence of F1 animals, as indicated by the PCA, was corroborated by the breed composition analysis, where the contribution of the Landrace and Yorkshire breeds showed that the F1 crossbreds were indeed hybrids of these two pure breeds. The variation in breed composition within the LYD population was not substantial, with standard deviations of 0.13, 0.12, and 0.19 for Landrace, Yorkshire, and Duroc, respectively. Similarly, the DK and LK populations exhibited balanced breed compositions. In the DK population, the average breed composition was 63% Duroc and 37% KNP, with minimal variation between individuals (SD = 0.05 for both breeds). The LK population had an average composition of 61% Landrace and 39% KNP, and low variation was also observed across individuals (SD = 0.06 for both breeds). These results suggest that the parental breeds had relatively balanced genetic contributions, as evidenced by the minimal variation in breed composition between individuals within the DK and LK populations.

**Genetic Parameter Estimates**

193    Heritability estimates for backfat thickness and carcass weight were derived from five different models; the

194    associated variance components are detailed in **Table 2**. The estimates of genetic additive variance ($V_g$) and error

195    variance ($V_e$) were used to calculate heritability for each trait.

196    Model 1 (NULL), which did not account for population structure, yielded the highest heritability estimates,

197    with a heritability value of $0.44 \pm 0.03$ for backfat thickness and $0.31 \pm 0.03$ for carcass weight. The elevated

198    heritability estimates for this model may be attributed to its lack of adjustments for potential confounding factors

199    related to breed differences. Models 2 (PCA_F) and 3 (GBC_F), which incorporated population structure as a

200    fixed effect, yielded lower heritability estimates; Model 2 estimated heritability for backfat thickness at $0.41 \pm$

201    $0.03$ and carcass weight at $0.26 \pm 0.03$, whereas Model 3 estimated these factors at $0.44 \pm 0.03$ and $0.27 \pm 0.03$,

202    respectively. These reductions in heritability suggest that accounting for population structure as a fixed effect can

203    decrease the perceived genetic influence on the traits. Models 4 (PCA_R) and 5 (GBC_R) included additional

204    genetic variance components ($V_{pc}$ and $V_{gbc}$) to account for population structure as a random effect. In Model 4,

205    the genetic variance ($V_g$) was estimated at $13.2 \pm 1.3$ and $V_{pc}$ at $1.6 \pm 1.7$ for backfat thickness, contributing an

206    additional heritability of $0.05 \pm 0.05$ to the base estimate of $0.41 \pm 0.04$. For carcass weight, $V_g$ was estimated at

207    $28.1 \pm 3.6$ and $V_{pc}$ at $23.7 \pm 15.1$, contributing an additional heritability of $0.19 \pm 0.1$ to the base estimate of $0.23$

208    $\pm 0.04$. Model 5 demonstrated similar patterns, although $V_{gbc}$ for backfat thickness was close to zero. These

209    models typically yielded heritability estimates similar to those of Model 1 for backfat thickness; however, for

210    carcass weight, they provided a more nuanced understanding of genetic effects by accounting for population

211    structure as a separate effect.

212

213    **Accuracy of Genomic Estimated Breeding Values**

214    The accuracy of GEBVs was evaluated using five models; the results are summarized in **Table 3** and depicted

215    in **Figure 3**. Model 1 (NULL), Model 4 (PCA_R), and Model 5 (GBC_R) exhibited the highest accuracy for

216    predicting both backfat thickness and carcass weight. These models achieved an average accuracy of 0.59 for

217    backfat thickness and 0.50 for carcass weight, with minimal variation across replicates (SD = 0.01 for backfat

218    thickness and between 0.03 to 0.04 for carcass weight).

219    Models that incorporated population structure as a fixed effect (Models 2 and 3) demonstrated lower accuracies

220    for GEBVs. For backfat thickness, Model 2 (PCA_F) achieved a mean accuracy of $0.40 \pm 0.03$, whereas Model 3

221    (GBC_F) yielded a mean accuracy of $0.53 \pm 0.04$. The accuracy for carcass weight in these models was reduced

222    similarly, with Model 2 achieving an accuracy of $0.34 \pm 0.03$ and Model 3 yielding an accuracy of $0.38 \pm 0.02$.

223    These results suggest that modeling population structure as a fixed effect captures population differences but

224    compromises GEBV accuracy. In contrast, modeling population structure as a random effect captures genetic

225    variation due to breed differences without adversely affecting GEBV accuracy.

226    The Spearman rank correlation coefficient of GEBV between all models showed that all models were highly

227    correlated with each other (except Model 2 in backfat thickness), ranging from 0.59 to 0.60. In carcass weight,

228    Models 1, 4, and 5 had high Spearman correlation coefficients with each other, but models 2 and 3 had low

229    correlation coefficients with the other models, ranging from 0.39 to 0.70 (**Figure 4**). Models that did not correct

230    for population structure and models that corrected for population structure as a random effect had similar genomic

231    prediction patterns.

232

233    # Discussion

234    In multi-breed genomic predictions, using a reference population that encompasses multiple breeds inevitably

235    introduces differences in population structure across these breeds. Therefore, this study aimed to assess prediction

236    accuracy while adjusting population structure as either a fixed or random effect in multi-breed genomic predictions.

237    The findings revealed that adjusting for population structure as a fixed effect resulted in decreased accuracy,

238    whereas treating it as a random effect did not yield any improvements in accuracy. These results suggest that in

239    multi-breed genomic predictions, the genomic relationship matrix sufficiently accounts for population structure,

240    indicating that a model without adjustments for population structure is the most efficient.

241

242    **Genotypic versus pedigree-based breed composition**

243    GBC highlights the superior accuracy of genotypic data over that of pedigree information in determining breed

244    composition. Pedigree records often contain inaccuracies or are incomplete, which can result in erroneous breed

245    composition estimates [21, 22]. In contrast, using genomic data with tools such as ADMIXTURE provides a more

246 precise assessment [23]. The findings of this study revealed that the breed compositions calculated using

247 ADMIXTURE closely aligned with those expected from complete pedigree records, thereby corroborating

248 previous research that emphasizes the reliability of genomic data for estimating breed composition in admixed

249 populations [23].

250

251 **Effect of population structure on genomic estimated breeding values**

252 The effect of population structure on the estimation of genetic parameters is a well-established concern in

253 genomic studies. Population structure can lead to false-positive associations [24], which may result in inflated

254 heritability estimates [10] and biased accuracies in genomic predictions [6]. To address this issue, this study

255 incorporated PCs and GBCs into GBLUP models as fixed or random effects.

256 Notably, the inclusion of PCs or GBCs as fixed effects resulted in decreased accuracy of GEBVs compared to

257 those of models that excluded these factors. This reduction in accuracy may stem from the redundancy between

258 the information provided by these variables and that captured by the genomic relationship matrix. Essentially, the

259 genomic relationship matrix already encompasses much of the population structure information; therefore, adding

260 PCs or breed composition as fixed effects could result in double-counting, leading to overcorrection and reduced

261 model accuracy [11, 25]. In contrast, treating PCs and GBC as random effects did not yield any improvement in

262 prediction accuracy. This result suggests that the additional genetic variance components captured by these

263 random effects did not provide significant new information beyond what was already accounted for by the

264 genomic relationship matrix. Similarly, previous studies have demonstrated that incorporating population

265 structure as a random effect does not enhance the accuracy of genomic predictions [25]. However, the advantage

266 of including breed as a random effect within the model, as GEBVs are divided into two components. Specifically,

267 a model with a random effect splits the genetic variance into within-breed and across-breed GEBVs, thereby

268 facilitating the understanding of how predictions differ within and across breeds [25].

269 These findings hold significant implications for the optimal design of genomic prediction models. Although

270 accounting for population structure is crucial to avoid biases, these results indicate that the genomic relationship

271 matrix within the GBLUP framework sufficiently captures the required information. Consequently, additional

272 adjustments for population structure, whether as fixed or random effects, may be unnecessary and could even

273 negatively affect prediction accuracy. These findings support the growing consensus that simpler models that rely

274 on the genomic relationship matrix without further correction for population structure are often the most effective

275 [25].

276     This study focused on carcass traits and therefore did not explicitly include heterozygosity, even though

277 crossbred animals were used. However, recent findings suggest that including heterozygosity in genomic

278 predictions for maternal traits can improve prediction accuracy [26]. Therefore, future research on maternal traits

279 in genomic prediction models may benefit from considering heterozygosity as a factor to further enhance

280 prediction accuracy.

281     **Implications for multi-breed genomic prediction**

282     Our findings have significant implications in the field of multi-breed genomic prediction. This study

283 demonstrated that the genomic relationship matrix alone could effectively capture breed differences within multi-

284 breed populations, thereby eliminating the necessity for additional corrections for population structure. This

285 circumvention is particularly advantageous in multi-breed contexts, where genetic relationships among breeds can

286 vary widely, facilitating accurate predictions of breeding values for selection decisions.

287     Given the observed decrease in accuracy when population structure was included as a fixed effect, future studies

288 and practical applications of genomic prediction should prioritize models that incorporate the genomic

289 relationship matrix as the primary tool for capturing genetic variance. This approach is more straightforward and

290 ensures higher accuracy in predicting breeding values, which is crucial for effectively managing and improving

291 crossbred populations.

292     In conclusion, this study underscores the robustness of the genomic relationship matrix in accounting for

293 population structure within multi-breed genomic prediction. The findings suggest that, although population

294 structure is an important consideration, the genomic relationship matrix is sufficient for capturing the relevant

295 genetic variance, modeling additional corrections unnecessary. This insight is valuable for optimizing genomic

296 prediction models in crossbred populations and enhancing the accuracy of GEBV predictions.

297

# References

298

299     1.    Goddard M. Genomic selection: prediction of accuracy and maximisation of long term response. Genetica.
300           2009;136(2):245-57.

301     2.    Thomasen JR, Guldbrandtsen B, Su G, Brøndum RF, Lund MS. Reliabilities of genomic estimated breeding
302           values in Danish Jersey. Animal. 2012;6(5):789-96.

303     3.    Olson KM, VanRaden PM, Tooker ME. Multibreed genomic evaluations using purebred Holsteins, Jerseys,
304           and Brown Swiss. J Dairy Sci. 2012;95(9):5378-83.

305     4.    Erbe M, Hayes B, Matukumalli L, Goswami S, Bowman P, Reich C, et al. Improving accuracy of genomic
306           predictions within and between dairy cattle breeds with imputed high-density single nucleotide
307           polymorphism panels. J Dairy Sci. 2012;95(7):4114-29.

308     5.    Hoze C, Fritz S, Phocas F, Boichard D, Ducrocq V, Croiseau P. Efficiency of multi-breed genomic selection
309           for dairy cattle breeds with different sizes of reference population. J Dairy Sci. 2014;97(6):3918-29.

310     6.    Makowsky R, Pajewski NM, Klimentidis YC, Vazquez AI, Duarte CW, Allison DB, et al. Beyond Missing
311           Heritability: Prediction of Complex Traits. Plos Genet. 2011;7(4).

312     7.    Price AL, Patterson NJ, Plenge RM, Weinblatt ME, Shadick NA, Reich D. Principal components analysis
313           corrects for stratification in genome-wide association studies. Nature Genetics. 2006;38(8):904-9.

314     8.    Pritchard JK, Donnelly P. Case-control studies of association in structured or admixed populations. Theor
315           Popul Biol. 2001;60(3):227-37.

316     9.    Marchini J, Cardon LR, Phillips MS, Donnelly P. The effects of human population structure on large genetic
317           association studies. Nature Genetics. 2004;36(5):512-7.

318     10.   Visscher PM, Yang JA, Goddard ME. A Commentary on 'Common SNPs Explain a Large Proportion of the
319           Heritability for Human Height' by Yang et al. (2010). Twin Res Hum Genet. 2010;13(6):517-24.

320     11.   Janss L, de Los Campos G, Sheehan N, Sorensen D. Inferences from genomic models in stratified populations.
321           Genetics. 2012;192(2):693-704.

322     12.   Burgos-Paz W, Souza CA, Megens HJ, Ramayo-Caldas Y, Melo M, Lemús-Flores C, et al. Porcine
323           colonization of the Americas: a 60k SNP story. Heredity. 2013;110(4):321-30.

324     13.   Park J, Kim Y, Jung H-J, Park B, Lee J, Moon H. Comparison of meat quality and physicochemical
325           characteristics of pork between Korean native black pigs (KNBP) and Landrace by market weight. J Anim
326           Sci Technol. 2005;47(1):91-8.

327     14.   Patterson N, Price AL, Reich D. Population structure and eigenanalysis. Plos Genet. 2006;2(12):2074-93.

328   15. Alexander DH, Novembre J, Lange K. Fast model-based estimation of ancestry in unrelated individuals.
329       Genome Res. 2009;19(9):1655-64.

330   16. Shringarpure S, Xing EP. Effects of Sample Selection Bias on the Accuracy of Population Structure and
331       Ancestry Inference. G3-Genes Genom Genet. 2014;4(5):901-11.

332   17. Yang JA, Lee SH, Goddard ME, Visscher PM. GCTA: A Tool for Genome-wide Complex Trait Analysis. Am
333       J Hum Genet. 2011;88(1):76-82.

334   18. Scrucca L, Fraley C, Murphy TB, Raftery AE. Model-based clustering, classification, and density estimation
335       using mclust in R: Chapman and Hall/CRC; 2023.

336   19. Lee SH, Van der Werf JH. MTG2: an efficient algorithm for multivariate linear mixed model analysis based
337       on genomic information. Bioinformatics. 2016;32(9):1420-2.

338   20. Lourenco DAL, Fragomeni BO, Tsuruta S, Aguilar I, Zumbach B, Hawken RJ, et al. Accuracy of estimated
339       breeding values with genomic information on males, females, or both: an example on broiler chicken.
340       Genetics Selection Evolution. 2015;47.

341   21. Kuehn LA, Keele JW, Bennett GL, McDaneld TG, Smith TPL, Snelling WM, et al. Predicting breed
342       composition using breed frequencies of 50,000 markers from the US Meat Animal Research Center 2,000
343       Bull Project. J Anim Sci. 2011;89(6):1742-50.

344   22. Funkhouser SA, Bates RO, Ernst CW, Newcom D, Steibel JP. Estimation of genome-wide and locus-specific
345       breed composition in pigs. Transl Anim Sci. 2017;1(1):36-44.

346   23. Gobena M, Elzo MA, Mateescu RG. Population Structure and Genomic Breed Composition in an Angus-
347       Brahman Crossbred Cattle Population. Front Genet. 2018;9.

348   24. Price AL, Zaitlen NA, Reich D, Patterson N. New approaches to population stratification in genome-wide
349       association studies. Nat Rev Genet. 2010;11(7):459-63.

350   25. Hayes BJ, Copley J, Dodd E, Ross EM, Speight S, Fordyce G. Multi-breed genomic evaluation for tropical
351       beef cattle when no pedigree information is available. Genetics Selection Evolution. 2023;55(1):71.

352   26. Iversen MW, Nordbø Ø, Gjerlaug-Enger E, Grindflek E, Lopes MS, Meuwissen T. Effects of heterozygosity
353       on performance of purebred and crossbred pigs. Genetics Selection Evolution. 2019;51:1-13.

354

355

356 # **Tables**

357 **Table 1. Genomic breed composition by breeds.**

| Population | Breed | Min | Median | Max | Mean | SD |
|---|---|---|---|---|---|---|
| Landrace × Yorkshire × Duroc | Landrace | 0.06 | 0.27 | 0.90 | 0.31 | 0.13 |
| | Yorkshire | 0.05 | 0.30 | 0.89 | 0.33 | 0.12 |
| | Duroc | 0 | 0.43 | 0.75 | 0.36 | 0.19 |
| Duroc × KNP | Duroc | 0.49 | 0.63 | 0.75 | 0.63 | 0.05 |
| | KNP | 0.25 | 0.37 | 0.51 | 0.37 | 0.05 |
| Landrace × KNP | Landrace | 0.43 | 0.62 | 0.75 | 0.61 | 0.06 |
| | KNP | 0.25 | 0.38 | 0.57 | 0.39 | 0.06 |

358

359

**Table 2. Variance components and heritability estimates from five models for backfat thickness and carcass weight traits. Variance components are the genetic additive variance ($V_g$) and the error variance ($V_e$). In addition, the Model 4 (PC_R) and the Model 5 (GBC_R) estimates additional genetic variance components ($V_{pc}$ and $V_{gbc}$).**

| Model | | Variance components | | Heritabilities | |
|---|---|---|---|---|---|
| | | Backfat thickness (mm) | Carcass weight (kg) | Backfat thickness (mm) | Carcass weight (kg) |
| 1 (NULL) | $V_g$ | 13.5 ± 1.3 | 31.4 ± 3.6 | 0.44 ± 0.03 | 0.31 ± 0.03 |
| | $V_e$ | 17.1 ± 0.8 | 69.2 ± 2.7 | | |
| 2 (PC_F) | $V_g$ | 12.1 ± 1.3 | 24.9 ± 3.7 | 0.41 ± 0.03 | 0.26 ± 0.03 |
| | $V_e$ | 17.5 ± 0.8 | 71.3 ± 2.8 | | |
| 3 (GBC_F) | $V_g$ | 13.7 ± 1.3 | 26.0 ± 3.4 | 0.44 ± 0.03 | 0.27 ± 0.03 |
| | $V_e$ | 17.1 ± 0.8 | 70.6 ± 2.7 | | |
| 4 (PC_R) | $V_g$ | 13.2 ± 1.3 | 28.1 ± 3.6 | 0.41 ± 0.04 | 0.23 ± 0.04 |
| | $V_{pc}$ | 1.6 ± 1.7 | 23.7 ± 15.1 | 0.05 ± 0.05 | 0.19 ± 0.1 |
| | $V_e$ | 17.2 ± 0.8 | 69.9 ± 2.7 | | |
| 5 (GBC_R) | $V_g$ | 13.5 ± 1.3 | 27.1 ± 3.5 | 0.44 ± 0.03 | 0.23 ± 0.04 |
| | $V_{gbc}$ | 0.2 ± 0.3 | 22.3 ± 14.3 | 0 | 0.19 ± 0.1 |
| | $V_e$ | 17.1 ± 0.8 | 70.2 ± 2.8 | | |

366 **Table 3. Mean and standard deviation of GEBV accuracy for five prediction methods.**

| Model | Backfat thickness (mm) | | Carcass weight (kg) | |
|---|---|---|---|---|
| | Mean | SD | Mean | SD |
| 1 (NULL) | 0.59 | 0.01 | 0.50 | 0.04 |
| 2 (PCA_F) | 0.40 | 0.03 | 0.34 | 0.03 |
| 3 (GBC_F) | 0.53 | 0.04 | 0.38 | 0.02 |
| 4 (PCA_R) | 0.59 | 0.01 | 0.50 | 0.03 |
| 5 (GBC_R) | 0.59 | 0.01 | 0.50 | 0.03 |

367

368

369

370 **Figures**

371



372

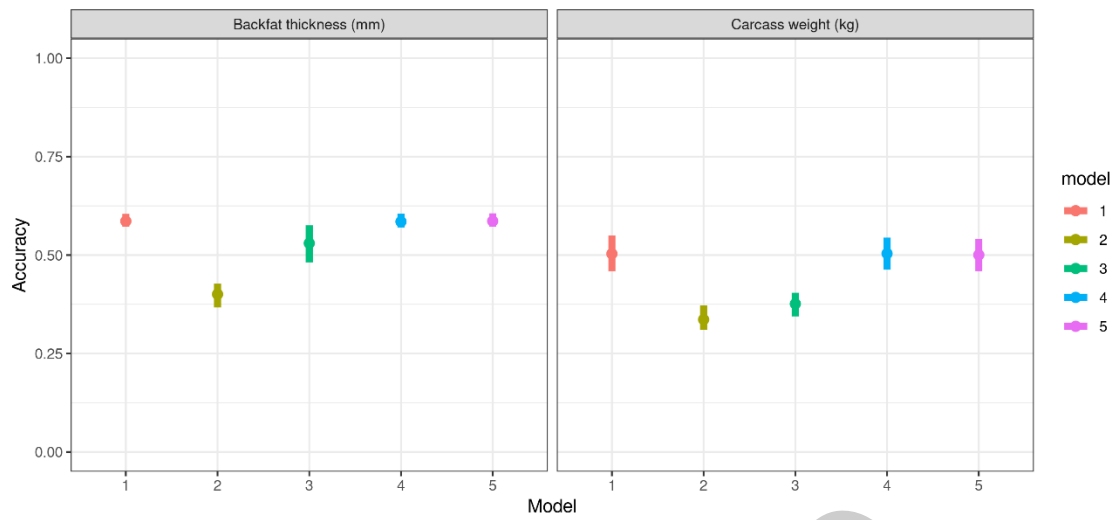**Figure 1. Population distribution across the first and second principal components.**

374

375

**Figure 2. Bar plot of the $Q$ matrix from an ADMIXTURE run, showing the proportion of the genome contributed by each breed. A shows the LYD population, B shows the DK population, and C shows the LK population. Each vertical bar represents an individual.**
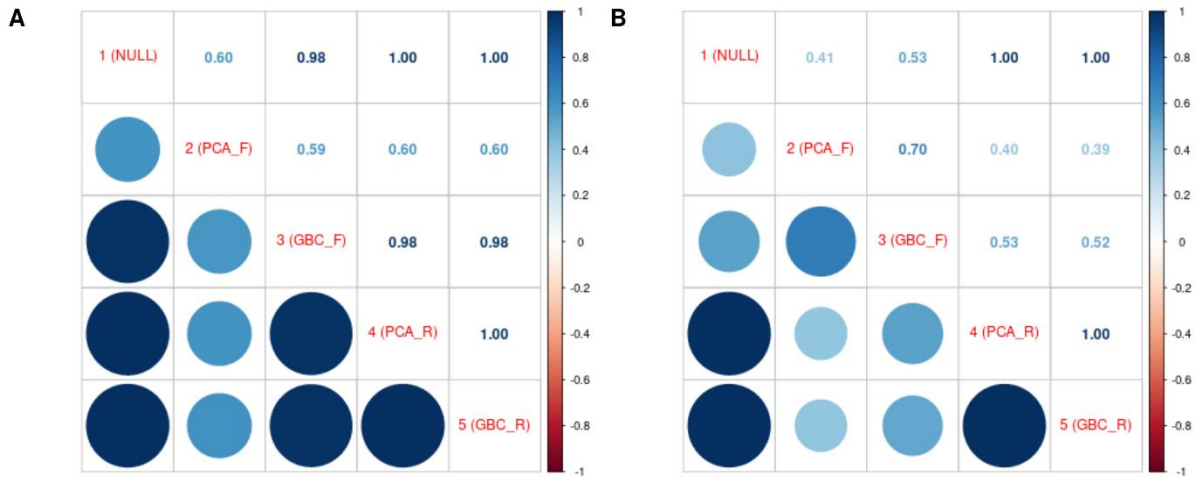
379

380

**Figure 3. GEBV accuracy of five prediction models. From left to right, the models are Model 1 (NULL), Model 2 (PCA_F), Model 3 (GBC_F), Model 4 (PCA_R), and Model 5 (GBC_R). The dots represent the average accuracy, and the lines indicate the standard deviation.**

384

385

**Figure 4. Spearman correlation between models. A represents backfat thickness and B represents carcass weight.**

386

387

388

389

390

391