



GenTORE

Genomic management Tools to Optimise Resilience and Efficiency

Grant agreement n°: 727213

H2020 - Research and Innovation Action

D4.2

Paper on breed admixture in the multi-breed data

Due date: M42

Actual submission date: M42

Project start date: 1st June 2017 **Duration:** 60 months

Workpackage concerned: WP4

Concerned workpackage leader: Mogens Sandø Lund

Lead Beneficiary: AU

Dissemination level:

- PU:** Public (must be available on the website)
- CO:** Confidential, only for members of the consortium (including the Commission Services)
- CI:** Classified, as referred to in Commission Decision 2001/844/EC



Table of content

1. Summary	3
2. Introduction	3
3. Results.....	3
4. Conclusions	3
5. Partners involved in the work	3
6. Annexes.....	4

1. Summary

New genomic prediction methods were developed and tested with simulated data. Genotypic data from real dairy populations, i.e. Danish Holstein, Swedish Red and Danish Jersey, were used as base populations to start genotype simulations for pure breeds and an admixed population. Phenotypes were also simulated, considering two traits which differ in their heritability levels, and different levels of genetic correlations between the breeds.

2. Introduction

Crossbreeding is an efficient strategy in dairy cattle breeding, to achieve better productivity and robustness at the animal and system level. Systems relying on crossbreeding, e.g. ProCROSS system, results in crossbred animals being highly admixed in terms of their breed of origin. Moreover, genomic evaluations in dairy cattle are generally carried out separately for pure breeds, and neither crossbred individuals' data are used, nor they get evaluations. This is partially due lack of methods which can efficiently handle data from pure breeds and admixed individuals. This deliverable provides and discusses such methodologies, which can handle data from multiple pure breeds and admixed individuals, allowing simultaneous evaluation of pure breeds and crossbred animals. The developed methodology relies on accurate estimation of breed origin of each genome segment (represented with a single marker), whose use in genomic prediction models has been proposed earlier, and generally known as BOA (breed origin of alleles).

3. Results

Results for proposed methods using simulated data are presented in Annex, in the format that it was submitted to a peer-reviewed journal, where it is currently under review.

4. Conclusions

The use of admixed individuals' data together with pure breeds' data in genomic prediction has two main consequences: (i) it increases the data size for all pure breeds, particularly for the breed with a small population size, allowing more accurate estimation of breeding values, (ii) it increases the prediction accuracy for admixed individuals.

5. Partners involved in the work

AU, ALLICE (current work)

6. Annexes

Starts on next page.

RESEARCH

Genomic prediction using a reference population of multiple pure breeds and admixed individuals

Emre Karaman^{1*}, Guosheng Su¹, Iola Croue² and Mogens S Lund¹

*Correspondence: emre@qgg.au.dk

¹Center for Quantitative Genetics and Genomics, Aarhus University, 8830 Tjele, Denmark

Full list of author information is available at the end of the article

Abstract

Background: In dairy cattle populations where crossbreeding has been used, animals show some level of diversity in their origins. In rotational crossbreeding, for instance, crossbred dams are mated with purebred sires from different pure breeds, and the genetic composition of crossbred animals is an admixture of the breeds included in the rotation. How to use the data of such individuals in genomic evaluations is still an open question. In this study, we aimed at providing methodologies for the use of data from crossbred individuals with admixed genetic background together with data from multiple pure breeds, for the purpose of genomic evaluations for both purebred and crossbred animals. A three-breed rotational crossbreeding system was mimicked using simulations.

Results: For purebred populations, within-breed genomic predictions generally led to higher accuracies than those from multi-breed predictions using combined data of pure breeds. Adding admixed population's (MIX) data to the combined pure breed data as of a different breed led to higher accuracies. When prediction models were able to account for breed origin of alleles, accuracies were generally higher than those from combining all available data, dependent on the correlation of QTL effects between the breeds. Accuracies varied when using SNP effects from any of the pure breeds to predict the breeding values of MIX. Using the breed-specific SNP effects estimated separately in each pure breed, while accounting for breed origin of alleles for the selection candidates of MIX, accuracies were generally improved compared to using SNP effects from multiple pure breed's reference data, but dependent on the correlation of QTL effects between the breeds. Models able to accommodate MIX data with breed origin of alleles approach generally led to higher accuracies than those from the models without considering breed origin of alleles.

Conclusions: Combining all available data, pure breed's and admixed population's data, in a multi-breed reference population is beneficial for the estimation of breeding values for pure breeds with a small reference population. For MIX, such an approach can lead to higher accuracies than considering breed origin of alleles, and using breed-specific SNP effects estimated separately in each pure breed. Including MIX data in the reference population of multiple breeds considering the breed origin of alleles, accuracies can be further improved. Our findings are relevant for breeding programs where crossing is systematically applied, and also for populations involving different subpopulations where exchange of genetic material among those became a routine.

Keywords: admixed population; crossbreeding; breed of origin; multi-breed prediction; genomic selection

1 **Background**

2 Genomic evaluation facilitates the accurate selection of genetically superior individ-
3 uals as early as their DNA samples are obtained [1]. Genetic progress by selection
4 depends on the accuracy of prediction. For genomic prediction, it depends on the
5 proportion of genetic variance explained by genome-wide single nucleotide polymor-
6 phisms (SNPs), and the accuracy with which the effect of those SNPs is estimated
7 [2, 3]. Both factors are conditional on the linkage disequilibrium (LD) between SNPs
8 and quantitative trait loci (QTL) [1–3].

9 For an accurate genomic prediction, a large population of individuals with both
10 phenotypes and genotypes is needed, which may not be possible for all traits and/or
11 all breeds [4–6]. In such cases, remedies would be to use SNP effects from another
12 breed (a strategy known as across-breed prediction) with a large reference popu-
13 lation, or to add data from other breeds (multi-breed prediction) to improve the
14 accuracy of SNP effect estimates. However, the accuracy of across-breed predic-
15 tions were generally around zero, and combining data from multiple breeds did not
16 notably improve accuracies in empirical studies [7–10].

17 When multiple breeds are combined to form a reference population, predictions
18 rely on the SNP-QTL LD across breeds. However, LD may be different [11, 12],
19 or the phase of the SNP and QTL alleles may be reversed [13] among the breeds,
20 due to selection and genetic drift [9]. The QTL, or SNPs in high linkage disequilib-
21 rium with QTL can be integrated into marker panels for genomic prediction with
22 a multi-breed reference population [14] or performing across-breed predictions [15].
23 Although this may alleviate the issue of SNP-QTL LD being different in different
24 breeds, it includes an implicit assumption that QTL effects are the same across the
25 breeds. This may not be true if, for instance, QTL by genetic background interac-
26 tions exist [10, 11]. Hence, it may be more appropriate to assume that QTL, and
27 therefore SNP effects are different but correlated than that they are the same across
28 the breeds.

29 Crossbreeding emerges as an efficient strategy for dairy cattle breeding to achieve
30 better productivity and robustness at the animal as well as the system level. The
31 improved performance is due to utilization of specific combining abilities and het-
32 erosis [16]. In dairy cattle populations, where crossbreeding has been used, animals
33 show different levels of diversity in their origins [11, 17]. In rotational crossbreed-
34 ing, for instance, where crossbred dams are mated to purebred sires from different
35 pure breeds, the genetic composition of crossbred animals is an admixture of the
36 breeds included in the rotation. At each rotation cycle, depending on the breed of
37 the sires used, admixture proportions of crossbred individuals changes greatly [18].
38 On the other hand, gene pool of some “purebred” populations may also contain
39 fractions from other breeds, because bulls are used across the breeds to some extent
40 [19]. A prerequisite for a well-structured crossbreeding system is to have an efficient
41 breeding plan within the pure breeds, as well as crossbred population. Because, a
42 sufficient number of purebred bulls are required for the system, and genetic gain in
43 the pure breeds should be maintained so as to ensure that overall economical benefit
44 over time is not negatively affected [20]. Nonetheless, genomic evaluations in dairy
45 cattle are mostly carried out separately for each breed, and neither cross breed data
46 is utilized nor breeders get genomic evaluations for their crossbred animals. It is,

47 therefore, required in some breeding programs that genomic prediction models are
48 able to accommodate a reference population including admixed individuals, as well
49 as multiple pure breeds, allowing simultaneous evaluation of all selection candidates.

50 An appealing approach to make use of data of admixed individuals in genomic
51 prediction is to incorporate breed proportions in genomic prediction models. Mak-
52 gahlela *et al.* [11] extended the random regression model to account for interactions
53 between marker effects and breed proportions, where the breed proportions were
54 inferred from pedigree in Nordic Red Dairy cattle. They reported that prediction
55 accuracy can be higher if breed proportions are considered. Thomassen *et al.* [21]
56 performed genomic predictions in Danish Jersey dairy cattle, and showed that a
57 model that accounts for breed proportions, estimated either from pedigree or mark-
58 ers, does not improve genomic predictions compared to a model that ignores it.
59 There are at least two limitations with both [11, 21] approaches. First, a single
60 measure of breed proportion may not be appropriate, because two individuals with
61 exactly the same breed proportions may have very different pattern of admixture
62 over their genome depending on which chromosomal region is inherited from which
63 pure breed [22]. Second, the correlations between the breeds were assumed to be
64 homogenous across the whole genome [21], or those correlations were even set to
65 zero due difficulties in estimation [11].

66 In this article, we propose a methodology suitable for genomic prediction using a
67 reference population of multiple purebred and admixed individuals. Using simula-
68 tions, we investigated the impact of correlation of QTL effects among the breeds,
69 and the heritability of the trait on accuracy of genomic prediction using different
70 approaches: (i) treating the combined data as of a single homogeneous population,
71 (ii) considering breed-specific SNP effects with/without accounting for correlations
72 between the breeds, and (iii) considering priors that lead to the use of region specific
73 correlations among the breeds.

74 **Methods**

75 **Data Simulation**

76 *Genotypes*

77 Genotype data at 51,477 loci were available for animals from each of the three
78 dairy cattle breeds: Danish Holstein (HOL), Swedish Red (RED) and Danish Jersey
79 (JER), from which a subset of 1,050 (HOL and RED) or 220 (JER) individuals were
80 formed the base populations for this study. The SNPs which were fixed for the same
81 allele in all three breeds were removed. For computational reasons, only the SNPs
82 (12,664) on first 5 chromosomes were considered. A plot summarising the principle
83 component analysis of genomic relationships among all animals was depicted to
84 assess the genetic relationships among the pure breeds (Additional File 1). In order
85 to establish a data set including multiple pure breeds (i.e., HOL, RED and JER)
86 and an admixed population (hereafter, MIX), a rotational crossbreeding system was
87 mimicked using simulations, considering three cycles of rotation (Table 1) for nine
88 generations. Using the same sets of base population genotype data, a total of 10
89 replicates were generated.

90 Simulations started with 1,050 (HOL and RED) or 220 (JER) individuals (gen-
91 eration 0 - G₀), of which 50 (HOL and RED) or 20 (JER) were assigned to be

92 males and the rest as females. The purebred populations were generated by mating
93 sires and dams from the same breed (Table 1). Population sizes and the number
94 of males and females were kept constant at each of the simulated generations for
95 HOL, RED and JER. This was achieved by mating 20 dams with the same sire, each
96 mating yielding one offspring, except for one mating which yielded two offsprings,
97 for the simulations of HOL and RED. In simulation of JER population, each sire
98 was mated with 10 dams, where each mating yielding one offspring, except for one
99 mating which yielded two offsprings.

100 The MIX in G1 was generated by mating sires from JER and dams from HOL of
101 G0. The MIX in G2 was generated by mating sires from RED and dams from MIX of
102 G1. Finally, one rotation cycle was completed with generating MIX in G3 by mating
103 sires from HOL and dams from MIX of G2. The following generations of MIX were
104 generated by mating sires from a pure breed, where the pure breed was dependent
105 on the rotation cycle, with the dams from the MIX (Table 1). Population size and
106 the number of males ($n = 50$) and females ($n = 1,050$) were also kept constant
107 at each of the simulated generations for MIX. When MIX individuals were mated
108 with HOL or RED, the mating structure was similar to that in those pure breeds,
109 whereas when MIX (or HOL) individuals were mated with JER, each JER sire
110 was mated with 50 dams, where 2 or 3 matings per sire was replicated to retain
111 the population size of MIX at 1,050. Selection was not considered, and mating was
112 completely at random.

113 The number of recombinations on each chromosome was determined using a ran-
114 dom variable drawn from a Poisson distribution, under the assumption that the
115 length of a chromosome in the Morgan's unit (we assumed $1 \text{ cM} \sim 1 \text{ Mb}$) is the
116 lambda parameter [23]. Recombination positions were sampled from a uniform dis-
117 tribution, and interference was ignored [23]. Mutation was not considered in the
118 simulations.

119 *Phenotypes*

120 The total number of QTL was set at 250, which were selected randomly among the
121 SNPs that satisfied $0.01 < \text{MAF} \leq 0.30$, where MAF is the minor allele frequency
122 computed as follows. First, allele frequency at each locus (p_k) was computed for each
123 breed, and then averaged over the breeds (\bar{p}_k), to avoid population sizes effecting
124 allele frequencies. Second, MAF of each locus was computed as $\min(\bar{p}_k, 1 - \bar{p}_k)$. The
125 selection of QTL with $0.01 < \text{MAF} \leq 0.30$ ensured that the QTL were segregating
126 with a lower MAF compared to SNPs, for the combined population at G0 (Table
127 2). The QTL were excluded from the final data set of SNPs. It should be noted
128 that although G0 was common to all 10 replicates, and therefore, the SNPs that
129 met the criteria to be selected as QTL were the same, the QTL or SNP sets did not
130 fully overlap among the replicates due to randomised selection of QTL. The effects
131 (explained below) of QTL were also simulated separately for each replicate.

132 Even if additive and dominance effects of QTL are the same in different breeds,
133 the difference in QTL allele frequencies may cause substitution effects of QTL [16]
134 to differ among the breeds, as well as genetic (co)variances. In this study, the QTL
135 effects were simulated directly from a multivariate normal distribution for varying
136 levels of correlations among the QTL effects of different breeds, that is correlations
137 of 1.00, 0.50 or 0.25.

138 Each individual had two alleles (maternal and paternal alleles) at each locus,
 139 inherited from its sire and dam. Breed origin of each allele for all loci were traced
 140 back to pure breeds at G0, and were known without error. Breeding value of each
 141 individual i (u_i) across G0-G9 were generated with:

$$u_i = \sum_{k=1}^{250} [Q_{ijk}^M * \alpha_{jk}^M + Q_{ijk}^P * \alpha_{jk}^P]$$

142 where Q_{ijk}^M and Q_{ijk}^P are the number of copies (0 or 1) of an arbitrarily chosen
 143 allele A at QTL locus k , inherited from its dam and sire breed j ($j=H,R,J$), respec-
 144 tively. The α_{jk}^M and α_{jk}^P are the simulated QTL effects for locus k , in breed j . The
 145 QTL effects were scaled such that the mean of the breed specific genetic variances
 146 (computed as the variance of breeding values) is 100 at G0. A random residual e_j
 147 drawn from a normal distribution, $e_j | \sigma_e^2 \sim N(0, \sigma_e^2)$, was added to each animal's
 148 breeding value to form phenotype. The size of σ_e^2 was determined according to the
 149 simulated heritabilities (explained later) and the mean genetic variance (100) over
 150 the breeds. It worths to note that, due to fixing the size of the residual variance
 151 across the breeds, heritabilities fluctuated around their mean values over the breeds.
 152 The same value of σ_e^2 was used in all generations for all individuals.

153 True (simulated) genetic correlations between the breeds were computed from
 154 the genetic variances, $\sigma_{u,j}^2 = \sum_{k=1}^{250} 2p_{jk}(1-p_{jk})\sigma_{\beta_j}^2$, and covariances, $\sigma_{u,jj'} =$
 155 $\sum_{k=1}^{250} \sqrt{2p_{jk}(1-p_{jk})2p_{j'k}(1-p_{j'k})}\sigma_{\beta_{jj'}}$ ($j=H,S,J$ and $j \neq j'$) [24] at k QTL.
 156 The genetic correlations between HOL-RED, HOL-JER and RED-JER were 0.88,
 157 0.75 and 0.78, respectively, for QTL correlation of 1.00, over 10 replicates and at
 158 G0. Those were 0.45, 0.38 and 0.38 for QTL correlation of 0.50, and 0.22, 0.19
 159 and 0.19 for QTL correlation of 0.25, respectively. The differences between QTL
 160 effect correlations and genetic correlations were due to the difference in QTL allele
 161 frequencies between the breeds. The correlations between QTL allele frequencies
 162 of HOL-RED, HOL-JER and RED-JER were 0.33, 0.22 and 0.41, respectively. The
 163 correlations between SNP allele frequencies were 0.47, 0.32 and 0.46. The QTL effect
 164 correlations of 0.50 and 0.25 are consistent with the reported genomic correlations
 165 (genetic correlations estimated based on available SNP set) among some cattle
 166 breeds for milk [14, 25] and fat [14], respectively. Two levels of heritability were
 167 considered for each scenario of correlations. Those were 0.40 and 0.05, which are of
 168 the same magnitude as the reported heritabilities of milk production and fertility
 169 traits, respectively (e.g.,[6]).

170 Reference and Validation Populations

171 Generations 6,7 and 8 (G6-G8) were used to form reference populations, while
 172 generation 9 (G9) was used to form validation populations. Hence, 660 JER indi-
 173 viduals, and 3,150 individuals from each of the HOL, RED and MIX were available
 174 in forming reference populations to estimate SNP effects.

175 The SNP effects were estimated using different reference populations: (i) a sin-
 176 gle pure breed (separate by breeds, i.e., **HOL**, **RED** or **JER**), (ii) combined data
 177 of multiple pure breeds (**HOL+RED+JER**), and (iii) combined data of multi-
 178 ple pure breeds and admixed (MIX) individuals. The MIX data was either used

179 as of a different “breed”, assuming homogeneous SNP effects across all breeds
 180 (**HOL+RED+JER+MIX**), or truly treated as an admixed population con-
 181 sidering breed origin of alleles (BOA) approach (**HOL+RED+JER+MIX un-**
 182 **cor/cor**).

183 The prediction of breeding values for each pure breed were performed using: (1)
 184 the estimated SNP effects from their own breed (within-breed prediction), (2) the
 185 estimated SNP effects from each of the other breeds (across-breed prediction), (3)
 186 the estimated SNP effects from a combined reference population (multi-breed pre-
 187 diction) and (4) the estimated SNP effects from a combined reference population
 188 considering BOA approach. The breeding values were predicted by multiplying SNP
 189 effects with allele dosages, with (4) or without (1-3) considering breed origin of
 190 alleles. These same strategies (1-4) were used to predict the breeding values of ad-
 191 mixed individuals. For the admixed individuals, SNP effects estimated separately
 192 using pure breed reference populations (**HOL/RED/JER**) were additionally used
 193 to predict breeding values, considering the BOA approach only for the validation
 194 animals (hereafter, pure-BOA). That is, breed origin of each SNP allele was traced
 195 back to its pure breed population only for the validation population, and the num-
 196 ber of counted alleles were multiplied by the breed specific estimate of SNP effects
 197 of the pure breeds.

198 We classified the methods using only a single breed’s data in model **training** to
 199 estimate SNP effects as **pure** (also includes pure-BOA as explained above), multiple
 200 breeds data without considering breed origin of alleles as **combined**, and multiple
 201 breed’s and MIX data considering breed origin of alleles as **BOA**.

202 Statistical Models

203 *Pure and combined*

204 A simple approach for genomic prediction using a combined reference population
 205 of multiple pure breeds and/or admixed individuals is to assume that the marker
 206 effects are the same across the breeds [26]. For this simple approach, when the
 207 data consisted of multiple breeds treated as of a single homogeneous population
 208 (Combined), we used the following model:

$$\mathbf{y} = \mathbf{1}\mu + \mathbf{X}\mathbf{b} + \mathbf{M}\boldsymbol{\beta} + \mathbf{e} \quad (1)$$

209 In the above equation, \mathbf{y} is a vector of phenotypes ($n \times 1$), $\mathbf{1}$ is the vector of 1s, μ
 210 is the general mean, \mathbf{X} is the matrix of breed proportions ($n \times 3$), \mathbf{b} is the vector
 211 of fixed breed effects (3×1), \mathbf{M} is the matrix of centered genotypes ($n \times l$) where
 212 centering was based on the current allele frequencies in the combined data, $\boldsymbol{\beta}$ is
 213 the vector of SNP effects, and \mathbf{e} is the vector of residuals ($n \times 1$). The value of n
 214 depends on the reference population size, and l is the number of SNPs. Model (1)
 215 was used without the breed proportions component $\mathbf{X}\mathbf{b}$ when the SNP effects were
 216 estimated separately for each pure breed (Pure and Pure-BOA).

217 *BOA*

218 Admixed breeds’ data was utilised by extending the existing linear model proposed
 219 for simple 2-way crosses (e.g., [27]) to accommodate more than two pure breeds:

$$\mathbf{y} = \mathbf{1}\mu + \mathbf{M}_H\boldsymbol{\beta}_H + \mathbf{M}_R\boldsymbol{\beta}_R + \mathbf{M}_J\boldsymbol{\beta}_J + \mathbf{e} \quad (2)$$

220 where \mathbf{y} is the vector of phenotypes ($n \times 1$) of all animals, that is, both purebred
 221 and admixed animals. The $\mathbf{1}$ is a vector of 1s, μ is the general mean, \mathbf{M}_H , \mathbf{M}_R and
 222 \mathbf{M}_J are the matrices of breed specific allele content of SNPs ($n \times l$) for HOL, RED
 223 and JER, respectively. The entry at a locus in, for instance \mathbf{M}_H , for an animal
 224 were the number (0,1 or 2) of counted allele A originated from HOL. That is, when
 225 the animal had no allele originating from HOL, or when a HOL animal had an aa
 226 genotype, the corresponding entry was zero. The same applied to matrices \mathbf{M}_R and
 227 \mathbf{M}_J . The matrices were column centered prior to analysis. The $\boldsymbol{\beta}_H$, $\boldsymbol{\beta}_R$ and $\boldsymbol{\beta}_J$ are
 228 vectors of SNP effects for HOL, RED and JER, respectively, and \mathbf{e} is the vector of
 229 residuals.

230 Bayesian Analysis

231 Bayesian approach was considered in parameter estimation, which requires assign-
 232 ing prior distributions to the unknowns of the model. Analyses were carried out
 233 separately for each trait. To investigate the impact of assuming heterogeneous
 234 (co)variance of SNP effects among different genome regions, three region sizes were
 235 considered based on a fixed number of SNPs; 1 SNP, 100 SNPs and the whole
 236 genome (WG). Regions sizes of 1 SNP and WG can be regarded as BayesA and
 237 SNP-BLUP [1] (or equivalently GBLUP [28]) when using model (1), and extensions
 238 of them for multiple components (breeds) when using model (2), respectively. In
 239 BayesA it is assumed that each SNP (1 SNP) follows a normal distribution with null
 240 mean and a locus-specific variance, while in GBLUP it is assumed that all SNPs
 241 (WG) have null means and a common variance. To consider heterogeneous variance
 242 of SNP effects among different genome regions using model (1), the matrix of geno-
 243 types and vector of SNP effects were partitioned into S subsets each with k_s loci
 244 ($s = 1, \dots, S$), and priors were assigned to each sub-vector of $\boldsymbol{\beta}$: $\boldsymbol{\beta}_s \mid \sigma_s^2 \sim N(\mathbf{0}, \mathbf{I}\sigma_s^2)$
 245 [29, 30]. The $\sigma_s^2(s)$ were further assigned a scaled inverse chi-square prior with a de-
 246 grees of freedom (df) and a scale parameter (S): $\sigma_s^2 \mid df, S \sim \chi^2(df, S)$. The values
 247 of hyper-parameters will be explained later.

248 In the analyses using model (2), all genotype matrices and vectors of SNP effects
 249 were also partitioned into S subsets each with l_s loci. A normal distribution prior
 250 was assigned for each sub-vector of SNP effects for population j ($j=H,R,J$): $\boldsymbol{\beta}_{j,s} \mid$
 251 $\sigma_{j,s}^2 \sim N(\mathbf{0}, \mathbf{I}\sigma_{j,s}^2)$. Hence, the SNP effects were breed-specific and uncorrelated
 252 across the breeds. That is, the genetic correlations among the breeds were assumed
 253 to be zero. The $\sigma_{j,s}^2(s)$ were further assigned a scaled inverse chi-square prior with a
 254 degrees of freedom (df_j) and a scale (S_j) parameter: $\sigma_{j,s}^2 \mid df_j, S_j \sim \chi^2(df_j, S_j)$. Using
 255 model (2), priors were also assigned such that the marker effects were breed-specific,
 256 but correlated between the breeds. That is, a multivariate normal distribution was
 257 assigned for each sub-vector of SNP effects: $[\boldsymbol{\beta}_{H,s} \boldsymbol{\beta}_{R,s} \boldsymbol{\beta}_{J,s}]' \mid \mathbf{B}_s \sim N(\mathbf{0}, \mathbf{B}_s \otimes \mathbf{I})$,
 258 where \mathbf{I} is an identity matrix of size equal to l_s if $l_s > 1$ or a scalar of 1 if $l_s = 1$.

$$\mathbf{B}_s = \begin{bmatrix} \sigma_{H,s}^2 & \sigma_{HR,s} & \sigma_{HJ,s} \\ \sigma_{RH,s} & \sigma_{R,s}^2 & \sigma_{RJ,s} \\ \sigma_{JH,s} & \sigma_{J,s} & \sigma_{J,s}^2 \end{bmatrix}$$

259 The diagonals of \mathbf{B}_s are the breed-specific SNP variances, while the off-diagonals
 260 are SNP covariances between the breeds. The \mathbf{B}_s was assumed to follow an inverted
 261 Wishart distribution with a shape (v_B) and a scale (\mathbf{V}_B) parameter: $\mathbf{B}_s \mid v_B, \mathbf{V}_B \sim$
 262 $IW(v_B, \mathbf{V}_B)$.

263 In both model (1) and (2), residuals were assigned a univariate normal prior,
 264 $e_i \mid \sigma_e^2 \sim N(0, \sigma_e^2)$, and variance σ_e^2 were assigned a scaled inverse chi-square prior
 265 with a degrees of freedom (df_e) and a scale (S_e) parameter: $\sigma_e^2 \mid df_e, S_e \sim \chi^2(df_e, S_e)$.
 266 Fixed effects were assigned flat priors.

267 The hyper-parameters of the prior distributions for the variance components were
 268 driven from the simulated genetic (co)variances and residual variances at G0
 269 as follows. For the analysis using model (2) assuming independent SNP effects among
 270 the breeds, $df_j = 4$ and $S_j = \frac{\sigma_{\beta_{j,s}}^2(df_j-2)}{df_j}$, where $\sigma_{\beta_{j,s}}^2 = \frac{\sigma_{u,j}^2}{\sum 2p_{j,l}(1-p_{j,l})}$ [31]. Here,
 271 $\sigma_{u,j}^2$ is the genetic variance for breed j , and p_{jl} is the allele frequency of l 'th SNP
 272 in breed j . Only one S_j was required for the analysis using model (1), which was
 273 computed using $\sigma_{u,j}^2$ (pure breed analysis) or the mean value of $\sigma_{u,j}^2$ over the breeds
 274 (combined analysis). For the analysis using model (2) assuming correlated SNP
 275 effects between the breeds, $\mathbf{V}_B = (v_B - 3 - 1)\mathbf{B}$ where $v_B = 6$, and

$$\mathbf{B} = \begin{bmatrix} \frac{\sigma_{u,H}^2}{\sum 2p_{H,j}(1-p_{H,j})} & \frac{\sigma_{u,HR}}{\sum \sqrt{2p_{H,j}(1-p_{H,j})}\sqrt{2p_{R,j}(1-p_{R,j})}} & \frac{\sigma_{u,HJ}}{\sum \sqrt{2p_{H,j}(1-p_{H,j})}\sqrt{2p_{J,j}(1-p_{J,j})}} \\ \frac{\sigma_{u,RH}}{\sum \sqrt{2p_{R,j}(1-p_{R,j})}\sqrt{2p_{H,j}(1-p_{H,j})}} & \frac{\sigma_{u,R}^2}{\sum 2p_{R,j}(1-p_{R,j})} & \frac{\sigma_{u,RJ}}{\sum \sqrt{2p_{R,j}(1-p_{R,j})}\sqrt{2p_{J,j}(1-p_{J,j})}} \\ \frac{\sigma_{u,JH}}{\sum \sqrt{2p_{J,j}(1-p_{J,j})}\sqrt{2p_{H,j}(1-p_{H,j})}} & \frac{\sigma_{u,JR}}{\sum \sqrt{2p_{J,j}(1-p_{J,j})}\sqrt{2p_{R,j}(1-p_{R,j})}} & \frac{\sigma_{u,J}^2}{\sum 2p_{J,j}(1-p_{J,j})} \end{bmatrix}$$

276 In the above equation, $\sigma_{u,j}^2$ and $\sigma_{u,jj'}$ ($j=H,R,J$ and $j \neq j'$) are genetic variances
 277 and covariances respectively. For residual variances, $df_e = 4$ and $S_e = \frac{\sigma_e^2(df_e-2)}{df_e}$,
 278 where σ_e^2 is the residual variance at G0.
 279

280 Markov-chain Monte Carlo (MCMC) algorithm with Gibbs sampling method was
 281 used to obtain samples of each parameter from its full conditional posterior distri-
 282 bution. Chain length for the analyses consisted of 50,000 cycles, of which the first
 283 10,000 were discarded as burn-in. Every 10th sample of the post burn-in cycles were
 284 stored for posterior analysis, yielding 4,000 posterior samples. Mean value of the
 285 posterior samples was used as the estimate of each parameter. All analysis were
 286 performed using self-written scripts in Julia [32].

287 Prediction accuracy

288 Prediction accuracy was assessed as the correlation between true and predicted
 289 breeding values of validation individuals (1,050 individuals for HOL, RED, MIX,
 290 and 220 individuals for JER) at G9. Accuracy of prediction using different data sets
 291 and models to estimate SNP effects were compared for each trait, QTL correlation
 292 and region size separately. Accuracy of prediction for different region sizes were
 293 compared for each data set and model, trait and QTL correlation, separately. All
 294 comparisons were performed using a two-sided paired t-tests, for which accuracies
 295 were paired across each replicate for the same validation population. A Bonferroni
 296 correction was used to control the Type 1 error rate of 0.05.

Results

Accuracies for all scenarios and all region sizes are given in Supplementary Tables 1-4, Additional file 2. For readability, only the core results from QTL effect correlation of 0.5 are presented in the main text. Accuracies were higher for high heritability trait than those for low heritability trait (Figures 1 and 2). Within-breed predictions in breeds with large reference populations (HOL and RED) were more accurate than in a breed with small reference population (JER). For the high heritability trait, within-breed predictions for HOL, RED and JER were 0.785, 0.747 and 0.629, respectively, when the region size was 1 SNP (Figure 1). For this high heritability trait, combining data from multiple pure breeds (HOL+RED+JER) assuming homogenous SNP effects (multi-breed prediction) did not improve, or even decreased (though not always significant) the accuracies for all breeds. Including the admixed population's (MIX) data in multi-breed prediction, as if it belongs to a different breed (HOL+RED+JER+MIX), yielded higher accuracies compared with combining only the data from pure breeds, and similar to or higher accuracies than using the single breed data alone (within-breed prediction), for genomic prediction of JER. When prediction models were able to accommodate data of admixed individuals by accounting for breed origin of alleles (HOL+RED+JER+MIX uncor/cor), accuracies were generally improved compared to combining all available data, but dependent on the correlation scenario. Across-breed predictions yielded accuracies much lower than within-breed predictions.

Accuracies were the lowest when using SNP effects from any of the pure breeds to predict the breeding values of admixed individuals. For the high heritability trait, predictions using SNP effects of HOL, RED and JER yielded accuracies of 0.411, 0.275 and 0.114, respectively, when the region size was 1 SNP (Figure 1). For the same scenario and region size, estimating SNP effects separately for each breed, but accounting for breed origin of alleles in prediction of breeding values (HOL/RED/JER) of MIX, improved accuracy up to 0.531. Combining MIX data with pure breeds' data assuming common SNP effects for all breeds (HOL+RED+JER+MIX), improved accuracies over combining only three pure breeds' data (HOL+RED+JER) for the accuracy of admixed individuals (0.792 vs 0.501). Models able to use MIX data with breed origin of alleles (HOL+RED+JER+MIX uncor/cor), improved accuracies over combining all available data, i.e., combining all purebred data or all purebred and admixed individuals' data, though dependent on the correlation of QTL scenario. For QTL correlation of 1.0, and predictions in MIX, (HOL+RED+JER+MIX) led to higher accuracies than (HOL+RED+JER+MIX uncor). Accounting (0.877) or not (0.876) for correlations between the SNP effects of different pure breeds did not make any difference (Figure 1). Among different region sizes considered here, region size of whole genome generally yielded the lowest accuracies for pure breeds and admixed population (Figure 3).

The importance of the methods considering breed origin of alleles in model training became more apparent as the correlation of the true QTL effects between the breeds decreased (See Supplementary Tables 1-4, Additional File 2). For high heritability trait and purebred populations, accuracies for (HOL+RED+JER+MIX uncor/cor) were significantly higher than those for (HOL+RED+JER+MIX) in QTL

343 effect correlation of 0.25. For MIX population, (HOL+RED+JER+MIX uncor/cor)
344 yielded significantly higher accuracies than those for (HOL+RED+JER+MIX), for
345 QTL effect correlation of 0.25, and for both traits.

346 Discussion

347 Within- and across-breed predictions

348 The simple approach for avoiding the unfavourable impact of the difference in
349 marker effects among the different purebred populations is to carry out separate
350 evaluations for each of those pure breeds, as is the case for genomic evaluations
351 in many countries [19]. Such an approach, however, comes with the cost of a po-
352 tential loss of data information resource, and therefore, in the accuracy of SNP
353 effect estimation. This is a limitation for genetic improvement in populations with
354 a small reference population. In this study, accuracies from within breed predictions
355 were higher for HOL and RED, compared to JER. Although there could be other
356 reasons, one explanation is the low reference population size (660 vs 3,150) set for
357 JER. The accuracies for pure breeds differed between the two heritability levels for
358 any QTL effect correlation scenario, with the high heritability trait having higher
359 accuracies than the low heritability trait. The fact that genomic prediction accuracy
360 is higher with large reference populations and/or for a high heritability trait was
361 reported in many other studies [3, 5, 33–35]. It should be noted that the accuracies
362 for the same heritability level fluctuated slightly for different QTL effect correlation
363 scenarios, because QTL effects were simulated using different multivariate normal
364 distributions (the covariance matrices differed) for those scenarios.

365 Using SNP effects of one pure breed to predict the breeding values of individu-
366 als of other breeds (across-breed prediction), yielded much lower accuracies than
367 within-breed predictions. This was true even when the simulated QTL effects had
368 a correlation of one. This is in line with the study of Steyn et al. [36] where several
369 breeds were simulated assuming identical QTL effects, but across-breed predictions
370 were poor. Studies using real data also showed that using data from one breed to
371 predict breeding values in other breeds results in accuracies as low as zero (e.g.,
372 [9, 10]). The prediction accuracy of MIX generally reflected the expected breed pro-
373 portions of the validation individuals. Using SNP effects from HOL, for instance,
374 led to the highest prediction accuracies for MIX, as HOL was the most recent an-
375 cestor population for MIX, and therefore, MIX individuals had a higher proportion
376 of their genome from HOL.

377 For within-breed predictions, both family relationships and linkage disequilibrium
378 (LD) between SNPs and QTL contribute to accuracy [37–39]. For across-breed pre-
379 diction, the relationships of the individuals of the target breed with the individuals
380 in the reference population are lower than those with the members of the target
381 breed. The relative contributions of the two factors, family relationships and LD,
382 to accuracy of breeding value estimation were not studied as it was not in the scope
383 of this paper. If we rely on the argument that low across-breed prediction accuracy
384 is due to differences in LD patterns among the breeds, that is the differences in
385 LD or the phase of the SNP and QTL alleles, then across-breed prediction can not
386 compete with within-breed prediction, even for closely related breeds. In addition
387 to LD patterns, it is also possible that QTL effects and/or QTL allele frequencies

388 differ among the breeds, while some QTL may only segregate in one breed [25, 40].
389 Needless to say, even if the QTL properties were the same among the breeds, SNP
390 effects would still be different to the extent to which LD between SNPs and QTL
391 differs between them [11, 12, 30, 41].

392 Although the simulated traits in this study were relatively polygenic, the variance
393 structure at the SNP level may be different from that of at the QTL level across
394 the genome [42, 43], favouring models that can accommodate such heterogeneity
395 [30, 44, 45]. The SNP panels tend to include SNPs with high minor allele frequency
396 (MAF), while the QTL have generally low MAF [46, 47]. The LD between the
397 two sets, SNPs and QTL, can not be perfect if their MAF differs. Because the
398 SNPs in a genome region are likely inherited together, and also likely to be in LD
399 with the same QTL, they may collectively capture the genetic variance at the QTL
400 [29, 45, 48]. Hence, assuming a common variance for groups of adjacent SNPs is
401 reasonable, while it allows more Bayesian learning compared to assuming variance
402 specific to every single SNP [49]. For region sizes larger than an optimum level, on
403 the other hand, the advantage of grouping adjacent SNPs will start to disappear as
404 the assumption on (co)variance will approach to that of whole genome region size
405 (WG).

406 For high heritability trait and purebred analysis, accuracies from different region
407 sizes were generally ranked as 100 SNPs > 1 SNP > WG. It was shown earlier by
408 simulations [30, 35] and real data analysis that assigning priors to groups of SNPs
409 may improve accuracies [44, 45] compared to assigning a common prior for all SNPs.
410 In a recent study, on the other hand, Liu *et al.*, [50] reported negligible differences
411 between several region sizes, 1- 30- or 100 SNPs and WG, for milk production and
412 fertility traits in Danish Jersey, and using a model which is nearly identical to our
413 model (1).

414 Combined data of multiple pure breeds

415 If the studied population is small, it might be challenging to establish a large ref-
416 erence population, and in turn the accuracy of genomic prediction might also be
417 limited [6]. For breeds with a limited reference population size, incorporating data
418 from other breeds may yield higher accuracies [26, 40, 51], though dependent on
419 the relatedness of those breeds [9, 19]. When HOL and RED individuals were in-
420 cluded in the reference population of JER (HOL+RED+JER), accuracies generally
421 dropped. Similarly, using that combined reference population, accuracies for HOL
422 and RED also generally dropped, but less compared to those for JER. When multi-
423 ple purebred populations are combined to form a reference population, SNP effects
424 are dominated by the breeds contributing more to the reference population. This
425 may cause prediction models to pick up only the effect of SNPs that are in LD with
426 QTL in all breeds, and/or only in the largest population, but not the effect of SNPs
427 specific to small populations [14]. We had additional simulations where all breeds
428 had the same number of individuals in the reference population (3,150 for each),
429 which led accuracies for JER also to be high and get less affected from the joint
430 analysis, as HOL and RED (results not given). These imply that the proportion of
431 each single breed in a combined reference population of multiple breeds is impor-
432 tant to achieve a sufficient accuracy for each breed, particularly when the breeds

433 are genetically distant. This was more formally investigated in [40] using a high
434 density SNP chip ($\sim 600,000$ SNPs), where one of the two breeds (Holstein and
435 Jersey) that formed a joint reference population had varying sizes, 0, 100, 500 or
436 2,000 animals, while the size of the other breed kept constant at 2,000 animals. As
437 the number of individuals of a breed in the joint reference population decreased,
438 accuracies for the candidates of the same breed also decreased [40].

439 In a study based on real genotypes of imputed sequence variants (~ 1 million
440 SNPs), van den Berg *et al.* [52] simulated phenotypes for four dairy cattle breeds
441 using identical QTL effects. They reported generally higher accuracies for multi-
442 breed predictions, compared to within-breed predictions. In our scenario of QTL
443 effect correlation of 1.0, the difference in the accuracies from within- and multi-
444 breed predictions were smaller compared to other (lower) QTL effect correlation
445 scenarios. At long distances of genome, LD differs between species and also between
446 different cattle breeds, whereas it is relatively consistent at short distances [3]. The
447 standard SNP sets, such as the one used here, are not sufficient to include all
448 such SNPs that are in high LD with QTL across the breeds. Moreover, we selected
449 QTL such that they had relatively low MAF compared to SNPs, whereas QTL
450 were randomly selected in [52], which have an impact on LD between QTL and
451 SNPs. These may partially explain why multi-breed genomic predictions generally
452 had lower accuracies than within-breed predictions even when the simulated QTL
453 effects were identical, compared to the findings of [52].

454 For the analysis of data consisting of multiple breeds (or lines, populations), an
455 appealing strategy is to apply multi-trait methods where the same trait in dif-
456 ferent breeds is considered as different but correlated traits, e.g. [8, 25]. In those
457 applications of multi-breed genomic prediction, however, a homogeneous genomic
458 correlation was assumed across the genome, for pairs of breeds. Lehermeier *et al.*
459 [41] applied a multivariate modelling approach, which is flexible in that both marker
460 effects and their (co)variances are allowed to differ among multiple breeds, but still
461 assumes a homogenous correlation across the genome of breed pairs. Chen *et al.* [53]
462 proposed a method which allows the estimation of SNP effects specific to each breed
463 while accounting for heterogenous (co)variances across the genome. Their method,
464 however, applies a variable selection procedure aiming to pinpoint the SNPs that
465 have an effect in all breeds involved, leaving out the SNPs with effect only on one
466 or a subset of the breeds. It was further extended by Calus *et al.* [10] so as to
467 accommodate also the selection of SNPs that are breed-specific. Nevertheless, both
468 methods [10, 53], make limited use of the correlated information in the data, be-
469 cause, regardless of how the SNPs to be included in the model are selected, their
470 effects are estimated separately within each breed. Furthermore, all those multi-trait
471 approaches are pertained to situations where individuals can be assigned to certain
472 pure breeds, and are not able to accommodate data of individuals with admixed
473 genetic background.

474 Genomic prediction including data from admixed individuals

475 If a large number of commercial farm data for admixed populations becomes avail-
476 able, it can help to improve selection accuracy by expanding the data size for each
477 pure breed population. Such data can also allow to exploit heterosis due domi-
478 nance, which would not be possible with purebred data [12]. How to use those data

479 in genomic evaluations is still an open question. Naturally, all purebred and
480 admixed individual data can be combined together, when homogeneous SNP effects
481 are assumed.

482 Including the data of admixed population (hereafter, MIX) along with the data of
483 pure breeds in the reference population led to higher accuracies than the combined
484 reference population of pure breeds. The JER benefited relatively more from adding
485 MIX data. Because we mimicked a rotational crossing system, at each generation,
486 admixed population individuals were sired by a purebred individual. Consequently,
487 when an admixed female was mated with a purebred male, the offspring had an en-
488 tire paternal chromosome from a pure breed, and a maternal chromosome including
489 large chunks of (i) admixture of all breeds and (ii) the pure breed of the maternal
490 grand-sire. This means that, at each generation following G1, pure breeds were not
491 equally represented in the genome of admixed individuals. Consider a single ad-
492 mixed individual at generation 6. That individual has expected breed composition
493 for a maternal chromosome of roughly 28% JER, 16% HOL, and 56% RED, and
494 for a paternal chromosome of 100% HOL. Those proportions change to be 14%
495 JER, 58% HOL, and 28% RED for a maternal chromosome, and 100% JER for a
496 paternal chromosome at generation 7, and to 57% JER, 29% HOL, and 14% RED
497 for a maternal chromosome, and 100% RED for a paternal chromosome at gener-
498 ation 8. Because a full rotation cycle of three generations (G6-G8) was considered
499 when forming the reference populations, each pure breed was represented in the
500 MIX data almost equally. Thereby, the reference population size indeed increased
501 almost equally for all breeds by adding MIX data to the combined data of three
502 breeds, HOL+RED+JER+MIX. As one would expect, JER benefited more from
503 this increase in data size, as it is the breed with the smallest pure breed refer-
504 ence population. It should be noted that the validation individuals of RED had the
505 grand-sires which were also the sires of MIX at G8, and G8 was included in the
506 reference population. Hence, although the data size increased almost equally for
507 each breed, the information in the data may not be equally informative for all the
508 breeds.

509 More elaborative ways to include individuals with admixed genetic background
510 in the genomic evaluations, were proposed. Makgahlela *et al.* [11] fitted a multi-
511 trait random regression model to account for interactions between marker effects
512 and breed proportions, where the breed proportions were inferred from pedigree in
513 Nordic Red Dairy cattle. They reported, for some traits, higher prediction accura-
514 cies for the model accounting for breed proportions, than a GBLUP model treating
515 the data as of a single homogeneous population. Another example of admixture is
516 admixture due to different populations, instead of breeds. Danish Jersey dairy cat-
517 tle, for instance, includes animals with different proportions of their genome from
518 original Danish and US Jersey populations [19, 21]. Although both originate from a
519 single breed, they have been separated long ago, and the persistency of phase were
520 shown to differ between the two, particularly at long distances of loci [21]. Hence,
521 the accuracy of genomic prediction for Danish Jersey may not only be challenged
522 by the small reference population size, but also by its admixed population struc-
523 ture. In order to overcome the negative impact of admixed population structure
524 in Danish Jersey on genomic prediction accuracy, Thomasen *et al.* [21] applied a

525 set of random regression models that included proportions of population origin for
526 each animal. Contrariwise, they [21] did not find any strong evidence that a model
527 which accounts for proportions of population origin, estimated either from pedigree
528 or markers, is superior to a model which ignores it. A possible explanation could
529 be that admixture due to different breeds may be a more important problem than
530 admixture due to subpopulations of the same breed, in genomic prediction. Nev-
531 ertheless, there are at least two limitations with both [11, 21] approaches. First,
532 breed proportions of an individual were average values over their whole genome,
533 because they were computed based solely on pedigree or markers. This may not be
534 appropriate, as two individuals with exactly the same breed proportions may have
535 very different admixture patterns over their genome depending on which chromo-
536 somal region is inherited from which pure breed [21, 22]. Second, their models are
537 somewhat restricted in that the correlations between the breeds were assumed to
538 be homogenous across the whole genome [21], or those correlations were even set to
539 zero for difficulties in estimation [11]. When the breeds are in different SNP-QTL
540 LD, the (co)variances of SNP effects are expected to differ over the genome, and
541 across the breeds [11, 21, 22, 41].

542 Genomic prediction considering breed origin of alleles

543 Models accounting for breed origin of each SNP allele, rather than genome-wide
544 breed proportions estimated from pedigree or markers, have been proposed, and
545 were shown to improve genomic predictions for simple 2 or 3-way crosses. Those
546 studies applied either univariate whole genome regression models at the SNP level
547 ignoring that the SNP effects might be correlated between the pure-breeds [27,
548 54], or rather computationally demanding multi-trait genomic BLUP models with
549 “partial” relationship matrices at the individual level [22, 55, 56]. It was claimed
550 that considering genomic correlations between pure-breed populations had limited
551 relevance in models for predicting crossbred performance [22, 55, 56].

552 Our results did not show any clear evidence of the benefit of accounting for corre-
553 lations among the breeds when MIX data were used with BOA approach, even for
554 the breed (JER) with a small reference population, for which one would expect more
555 gain in accuracy compared with breeds with a large reference population (HOL and
556 RED). A possible explanation of unobserved benefit for JER could be due to this
557 breed being genetically distinct from HOL and RED [57], and therefore, the pattern
558 of SNP effects over the genome being different from HOL and RED. Additionally,
559 the information in the data may be weak to estimate correlations among the breeds.
560 The MIX data also increased the within-breed data size to some extent, which may
561 lower the importance of correlated information from other breeds [41]. For the sce-
562 nario of QTL effect correlation 1.00, analysis with HOL+RED+JER+MIX was
563 competitive with or even superior to analysis using BOA without accounting for
564 correlations between the breeds, particularly in predicting breeding values MIX.
565 This is may be due including MIX individuals in the reference population simply
566 increasing the data size in a joint analysis, whereas BOA with uncorrelated analysis
567 utilize only the information in a single breed.

568 The differences in LD pattern and phase persistency across different breeds [43]
569 may result in marker effects to be highly correlated for regions, where LD and SNP-
570 QTL phase is constant between the breeds [41]. Hence, we have anticipated that

571 correlations among the populations at the region level might improve the accuracy
572 of genomic predictions, even though the correlations at the whole genome level do
573 not. In this study, the differences in accuracies from 100 SNPs and 1 SNP region
574 sizes were generally negligible, whereas WG generally yielded the lowest accuracies.
575 It worths to note that, however, the fixed-length of 100 SNPs as region size was
576 arbitrarily chosen to give an insight on the impact of grouping SNPs in within-,
577 across- and multi-breed genomic prediction accuracy, and there may exist other
578 region sizes to yield higher accuracies than 100 SNPs. In analysis aiming to utilize
579 correlations between the breeds, such as analysis using BOA approach, the knowl-
580 edge of the LD patterns and persistence of phase among the breeds may be useful
581 in grouping SNPs.

582 van den Berg *et al.*, [14] showed that prediction of breeding values and genomic
583 correlations across populations can be more accurate if a carefully selected set of
584 causal variants or SNPs that are very close to causal variants from sequencing data
585 are used together with commercial SNP panels. Doing so may alleviate the issue of
586 SNP-QTL LD being different in different breeds. In a recent study, Liu *et al.* [6]
587 showed that integrating additional selected sequence variants to the standard 54K
588 SNP chip led to significant improvements of reliabilities for the genomic evaluation
589 of milk production traits in Danish Jersey. They reported that the benefits of using
590 selected sequence variants in genomic prediction for milk and protein remained
591 significant even in the scenario of the largest reference population consisting of
592 animals from Danish and US Jersey populations. In order to eliminate the impact
593 of LD differences among the breeds on the comparison of accuracy for using the
594 two BOA approaches (correlated and uncorrelated SNP effects), we ran additional
595 analyses for the QTL correlation scenario of 0.50 and the low heritability trait ($h^2 =$
596 0.05), using only the 250 QTL as SNPs and the region size of 1 SNP. Analyses using
597 BOA approach assuming correlated SNP effects were higher than those assuming
598 uncorrelated SNP effects between the breeds (Figure 4). In light of these, one can
599 argue that integrating selected sequence variants may be an efficient way of using
600 correlated information from the breeds, and that in this case taking into account
601 the correlation of SNP effects between breeds may allow for greater accuracy, in
602 genomic evaluations with data from multiple purebred and admixed individuals
603 using BOA approach.

604 Estimation of the breed composition of individuals with admixed genomic back-
605 ground is of relevance for genomic prediction, because if not accounted for it may
606 lead to spurious estimates of SNP effects [58]. In real life applications, pedigree
607 records and/or parentage validation can be used to distinguish purebred and ad-
608 mixed animals, but any error in the pedigree may lead to inaccurate consideration
609 of individuals as pure or admixed [18]. Nevertheless, genomic prediction should rely
610 on local ancestry (i.e., breed of origin) for each of the SNP alleles, rather than
611 a genome-wide (global) ancestry computed from pedigree or markers [59]. Meth-
612 ods exist to estimate local ancestry in a population of admixed individuals (e.g.,
613 [60]). In this simulation study, breed origin of admixed individuals were known
614 without error, but those could also be estimated from the data of purebred in-
615 dividuals. Due to mimicking a systematic crossing scheme in our simulations with
616 well-defined purebred individuals, such estimates are expected to be highly accurate

617 (Ana C. Guillenea, personal communication). For populations, where admixture is
618 more complex, however, one first needs to find the number of pure breeds in the
619 gene pool, and then to assign breed origin to each SNP allele for all animals in the
620 population. This may introduce another source of error, and the models requiring
621 breed origin of alleles, with or without accounting for correlations, may suffer from
622 such errors to the extent where simply combining all available data (multiple pure
623 and admixed breeds data) might become highly competitive. It was shown that a
624 higher number of animals would be required to distinguish closely related breeds
625 than to distinguish distantly related breeds [61], when the breed origin of an animal
626 is needed to be inferred from the genotypic data. To the best of our knowledge,
627 there is no information on the number of purebred animals required to correctly
628 assign breed-origin of alleles of the crossbred animals.

629 Genome scaling

630 Approximations for genomic prediction accuracy [3, 62] use the size of the reference
631 population (n_R), trait heritability (h^2), and the effective number of chromosomal
632 segments segregating in the population (M_e), where M_e is a function of the genome
633 length and the effective population size (N_e). Following those studies [3, 62], within-
634 breed prediction accuracy can be estimated with $\sqrt{h^2 n_R / (h^2 n_R + M_e)}$. In this study,
635 only the first five chromosomes were simulated, which is roughly a quarter of the
636 cattle genome. Those approximations suggest that, if we scale up the genome size
637 (and the number of QTL) to that of the whole genome, and the size of the reference
638 populations accordingly, our results will still hold, in within-breed predictions. For
639 across-breed prediction, Wientjes et al. [63], suggested to use of $r_g \sqrt{h^2 n_R / (h^2 n_R + M_e)}$,
640 where r_g is the genetic correlation between the breeds. They further suggested that
641 M_e values of 20,000 and 40,000 may be used when the populations are closely and
642 distantly related, respectively. On the other hand, combining different breeds to-
643 gether will increase N_e [64], and thereby M_e , requiring a larger reference population
644 size to compensate this increase in M_e , to avoid a reduction in accuracy [36, 52].
645 Models accounting for BOA, on the other hand, make use of single-breed data, while
646 taking advantage of an increase in n_R by using data from admixed individuals. The
647 BOA model with correlations further utilizes correlated information from other
648 breeds. It worths to note that those approximations assume a single homogenous
649 target (validation) population.

650 Conclusion

651 The aim of this simulation study was to provide a model allowing the inclusion
652 of data from individuals with admixed genetic background in genomic evaluations,
653 while accounting for the differences of marker effects for purebred populations in
654 the gene pool. Combining pure breeds' and admixed population's data, in a multi-
655 breed reference population was beneficial for the estimation of breeding values for
656 pure breeds with a small reference population. For the admixed population, com-
657 bining all available data (from purebred and admixed individuals) and realizing a
658 combined genomic evaluation led to higher accuracies than considering BOA for se-
659 lection candidates only, and using breed-specific SNP effects estimated separately in
660 each pure breed. Including admixed individual's data in the reference population of

661 multiple breeds considering the BOA approach, accuracies were further improved.
 662 Our findings are relevant for breeding programs where crossing is systematically
 663 applied (e.g., ProCROSS system, <http://www.procross.info>), and also for popula-
 664 tions involving different subpopulations where exchange of genetic materials among
 665 those became a routine (e.g., Nordic Red dairy cattle).

666 Competing interests

667 The authors declare that they have no competing interests.

668 Availability of data and material

669 The datasets used during the current study are available from the corresponding author on reasonable request.

670 Funding

671 This project has received funding from the European Union's Horizon 2020 research and innovation programme -
 672 GenTORE - under grant agreement No 727213.

673 Author's contributions

674 EK simulated the data, contributed to the formulation of the methods, implemented the methods and performed
 675 the analysis, and drafted the manuscript. GS co-supervised the study, contributed to the formulation of the
 676 methods, and revised the manuscript. IC contributed to the design of data simulation and discussion of the results.
 677 MSL conceived and supervised the study, contributed to the formulation of the methods and discussion of the
 678 results. All authors read and approved the final manuscript.

679 Author details

680 ¹Center for Quantitative Genetics and Genomics, Aarhus University, 8830 Tjele, Denmark. ²ALLICE, F-78350
 681 Jouy-en-Josas, France.

682 References

- 683 1. Meuwissen, T., Hayes, B., Goddard, M.: Prediction of total genetic value using genome-wide dense marker
 684 maps. *Genetics* **157**, 1819–1829 (2001)
- 685 2. Dekkers, J.: Prediction of response to marker assisted and genomic selection using selection index theory. *J*
 686 *Anim Breed Genet* **124**, 591–611 (2007)
- 687 3. Goddard, M.: Genomic selection: prediction of accuracy and maximization of long term response. *Genetica*
 688 **136**(245-257) (2009)
- 689 4. Hayes, B., Bowman, P., Chamberlain, A., Verbyla, K., Goddard, M.: Accuracy of genomic breeding values in
 690 multi-breed dairy cattle populations. *Genet Sel Evol*, 41–51 (2009)
- 691 5. Karaman, E., Cheng, H., Firat, M., Garrick, D., Fernando, R.: An upper bound for accuracy of prediction using
 692 GBLUP. *PLoS ONE* **11**, 0161054 (2016)
- 693 6. Liu, A., Lund, M., Boichard, D., Karaman, E., Fritz, S., Aamand, G., Nielsen, U., Wang, Y., Su, G.:
 694 Improvement of genomic prediction by integrating additional single nucleotide polymorphisms selected from
 695 imputed whole genome sequencing data. *Heredity* (2019)
- 696 7. Erbe, M., Hayes, B., Matukumalli, L., Goswami, S., Bowman, P., Reich, C., Mason, B., Goddard, M.:
 697 Improving accuracy of genomic predictions within and between dairy cattle breeds with imputed high-density
 698 single nucleotide polymorphism panels. *J Dairy Sci* **95**(7), 4114–29 (2012)
- 699 8. Olson, K., VanRaden, P., Tooker, M.: Multibreed genomic evaluations using purebred Holsteins, Jerseys, and
 700 Brown Swiss. *J Dairy Sci* **95**(9), 5378–5383 (2012)
- 701 9. Kachman, S., Spangler, M., Bennett, G., Hanford, K., Kuehn, L., Snelling, W., Thallman, R., Saatchi, M.,
 702 Garrick, D., Schnabel, R., Taylor, J., Pollak, E.: Comparison of molecular breeding values based on within- and
 703 across-breed training in beef cattle. *Genet Sel Evol* **45**(1), 30 (2013)
- 704 10. Calus, M., Goddard, M., Wientjes, Y., Bowman, P., Hayes, B.: Multibreed genomic prediction using multitrait
 705 genomic residual maximum likelihood and multitask Bayesian variable selection. *J Dairy Sci* **101**(5), 4279–4294
 706 (2018)
- 707 11. Makgahlela, M., Mantysaari, E., Strandén, I., Koivula, M., Nielsen, U., Sillanpää, M., Juga, J.: Across breed
 708 multi-trait random regression genomic predictions in the Nordic Red dairy cattle. *J Anim Breed Genet* **130**(1),
 709 10–19 (2013)
- 710 12. Veroneze, R., Bastiaansen, J., Knol, E., Guimaraes, S., Silva, F., Harlizius, B., Lopes, M.S., Lopes, P.: Linkage
 711 disequilibrium patterns and persistence of phase in purebred and crossbred pig (*Sus scrofa*) populations. *BMC*
 712 *Genet* **15**(1), 126 (2014)
- 713 13. de Roos, A., Hayes, B., Spelman, R., Goddard, M.: Linkage disequilibrium and persistence of phase in
 714 Holstein-Friesian, Jersey and Angus cattle. *Genetics* **179**(3), 1503–1512 (2008)
- 715 14. van den Berg, I., Boichard, D., Lund, M.: Sequence variants selected from a multi-breed GWAS can improve
 716 the reliability of genomic predictions in dairy cattle. *Genet Sel Evol* **48**(1), 83–83 (2016)
- 717 15. Raymond, B., Bouwman, A., Schrooten, C., Houwing-Duistermaat, J., Veerkamp, R.: Utility of whole-genome
 718 sequence data for across-breed genomic prediction. *Genet Sel Evol* **50**(1), 27 (2018)
- 719 16. Falconer, D., Mackay, T.: *Introduction to Quantitative Genetics*, 4th edn. Benjamin Cummings, ??? (1996)
- 720 17. Hess, M., Druet, T., Hess, A., Garrick, D.: Fixed-length haplotypes can improve genomic prediction accuracy in
 721 an admixed dairy cattle population. *Genet Sel Evol* **49**(1), 54 (2017)
- 722 18. Crum, T., Schnabel, R., Decker, J., Regitano, L., Taylor, J.: CRUMBLER: A tool for the prediction of ancestry
 723 in cattle. *PLoS ONE* **14**(8), 0221471 (2019)

- 724 19. Mogens, S., Guosheng, S., Luc, J., Bernt, G., Rasmus, F.: Genomic evaluation of cattle in a multi-breed
725 context. *Livest Sci* **166**, 101–110 (2014)
- 726 20. Sørensen, M., Norberg, E., Pedersen, J., Christensen, L.: Invited review: Crossbreeding in dairy cattle: A Danish
727 perspective. *J Dairy Sci* **91**(11), 4116–4128 (2008)
- 728 21. Thomsen, J., Sørensen, A., Su, G., Madsen, P., Lund, M., Guldbandsen, B.: The admixed population
729 structure in Danish Jersey dairy cattle challenges accurate genomic predictions. *J Anim Sci* **91**(7), 3105–3112
730 (2013)
- 731 22. Sevillano, C.A., ten Napel, J., Guimaraes, S., Silva, F., Calus, M.: Effects of alleles in crossbred pigs estimated
732 for genomic prediction depend on their breed-of-origin. *BMC Genomics* **19**(1), 740 (2018)
- 733 23. Yamamoto, E., Matsunaga, H., Onogi, A., Kajiya-Kanegae, H., Minamikawa, M., Suzuki, A., Shirasawa, K.,
734 Hirakawa, H., Nunome, T., Yamaguchi, H., Miyatake, K., Ohyama, A., Iwata, H., Fukuoka, H.: A
735 simulation-based breeding design that uses whole-genome prediction in tomato. *Sci Rep* **6**(19454) (2016)
- 736 24. Li, X., Lund, M., Janss, L., Wang, C., Ding, X., Zhang, Q., Su, G.: The patterns of genomic variances and
737 covariances across genome for milk production traits between Chinese and Nordic Holstein populations. *BMC*
738 *Genet* **18**(26), 12 (2017)
- 739 25. Zhou, L., Lund, M., Wang, Y., Su, G.: Genomic predictions across Nordic Holstein and Nordic Red using the
740 genomic best linear unbiased prediction model with different genomic relationship matrices. *J Anim Breed*
741 *Genet* **131**(4), 249–257 (2014)
- 742 26. de Roos, A., Hayes, B., Goddard, M.: Reliability of genomic predictions across multiple populations. *Genetics*
743 **183**, 1545–1553 (2009)
- 744 27. Ibánñez-Escriche, N., Fernando, R., Toosi, A., Dekkers, J.: Genomic selection of purebreds for crossbred
745 performance. *Genet Sel Evol* **41**(1), 12 (2009)
- 746 28. Strandén, I., Garrick, D.: Technical note: Derivation of equivalent computing algorithms for genomic predictions
747 and reliabilities of animal merit. *J Dairy Sci* **92**(6), 2971–2975 (2009)
- 748 29. Zeng, J., Garrick, D., Dekkers, J., Fernando, R.: A nested mixture model for genomic prediction using
749 whole-genome SNP genotypes. *PLoS ONE* **13**(3), 0194683 (2018)
- 750 30. Karaman, E., Lund, M., Su, G.: Multi-trait single-step genomic prediction accounting for heterogeneous
751 (co)variances over the genome. *Heredity* **124**, 274–287 (2020)
- 752 31. Habier, D., Fernando, R.L., Kizilkaya, K., Garrick, D.J.: Extension of the Bayesian alphabet for genomic
753 selection. *BMC Bioinform* **12**(1), 186 (2010)
- 754 32. Bezanson, J., Edelman, A., Karpinski, S., Shah, V.: Julia: A fresh approach to numerical computing. *SIAM*
755 *Review* **59**(1), 65–98 (2017)
- 756 33. Daetwyler, H., Pong-Wong, R., Villanueva, B., Woolliams, J.: The impact of genetic architecture on
757 genome-wide evaluation methods. *Genetics* **185**(3), 1021–1031 (2010)
- 758 34. Cheng, H., Kizilkaya, K., Zeng, J., Garrick, D., Fernando, R.: Genomic prediction from multiple-trait Bayesian
759 regression methods using mixture priors. *Genetics* **209**(1), 89–103 (2018)
- 760 35. Karaman, E., Lund, M., Anche, M., Janss, L., Su, G.: Genomic prediction using multi-trait weighted GBLUP
761 accounting for heterogeneous variances and covariances across the genome. *G3-Genes Genom Genet* **8**(11),
762 3549–3558 (2018)
- 763 36. Steyn, Y., Lourenco, D., Misztal, I.: Genomic predictions in purebreds with a multibreed genomic relationship
764 matrix. *J Anim Sci* **97**(11), 4418–4427 (2019)
- 765 37. Habier, D., Fernando, R., Dekkers, J.: The impact of genetic relationship information on genome-assisted
766 breeding values. *Genetics* **177**, 2389–2397 (2007)
- 767 38. Habier, D., Tetens, J., Seefried, F.-R., Lichtner, P., Thaller, G.: The impact of genetic relationship information
768 on genomic breeding values in German Holstein cattle. *Genet Sel Evol* **42**(1), 5 (2010)
- 769 39. Wientjes, Y., Veerkamp, R., Calus, M.: The effect of linkage disequilibrium and family relationships on the
770 reliability of genomic prediction **193**(2), 621–631 (2013)
- 771 40. Wientjes, Y., Calus, M., Goddard, M., Hayes, B.: Impact of QTL properties on the accuracy of multi-breed
772 genomic prediction. *Genet Sel Evol* **47**(42) (2015)
- 773 41. Lehermeier, C., Schön, C., de Los Campos, G.: Assessment of genetic heterogeneity in structured plant
774 populations using multivariate whole-genome regression models. *Genetics* **201**(1), 323–337 (2015)
- 775 42. de los Campos, G., Sorensen, D., Gianola, D.: Genomic heritability: What is it? *PLoS Genet* **11**, 1005048
776 (2015)
- 777 43. Wang, L., Sørensen, P., Janss, L., Ostensen, T., Edwards, D.: Genome-wide and local pattern of linkage
778 disequilibrium and persistence of phase for 3 Danish pig breeds. *BMC Genet* **14**(1), 115 (2013)
- 779 44. Brøndum, R.F., Su, G., Lund, M., Bowman, P., Goddard, M., Hayes, B.: Genome position specific priors for
780 genomic prediction. *BMC Genomics* **13**(1), 543 (2012)
- 781 45. Gebreyesus, G., Lund, M., Buitenhuis, B., Bovenhuis, H., Poulsen, N., Janss, L.: Modeling heterogeneous
782 (co)variances from adjacent-SNP groups improves genomic prediction for milk protein composition traits. *Genet*
783 *Sel Evol* **49**(1), 89 (2017)
- 784 46. Goddard, M., Hayes, B.: Mapping genes for complex traits in domestic animals and their use in breeding
785 programmes. *Nat Rev Genet* **10**, 381–391 (2009)
- 786 47. Kemper, K., Goddard, M.: Understanding and predicting complex traits: knowledge from cattle. *Hum Mol*
787 *Genet* **21**(1), 45–51 (2012)
- 788 48. Sørensen, L., Janss, L., Madsen, P., Mark, T., Lund, M.: Estimation of (co)variances for genomic regions of
789 flexible sizes: application to complex infectious udder diseases in dairy cattle. *Genet Sel Evol* **44**(1), 18 (2012)
- 790 49. Gianola, D., de los Campos, G., Hill, W., Manfredi, E., Fernando, R.: Additive genetic variability and Bayesian
791 alphabet. *Genetics* **183**(1), 347–363 (2009)
- 792 50. Liu, A., Lund, M., Boichard, D., Karaman, E., Guldbandsen, B., Fritz, S., Aamand, G., Nielsen, U., Sahana,
793 G., Wang, Y., Su, G.: Weighted single-step genomic best linear unbiased prediction integrating variants selected
794 from sequencing data by association and bioinformatics analyses. *Genet Sel Evol* **52**(1), 48 (2020)
- 795 51. Lund, M., de Roos, A., de Vries, A., Druet, T., Ducrocq, V., Fritz, S., Guillaume, F., Guldbandsen, B., Liu,

- 796 Z., Reents, R., Schrooten, C., Seefried, F., Su, G.: A common reference population from four European
797 Holstein populations increases reliability of genomic predictions. *Genet Sel Evol* **43**(1), 43 (2011)
- 798 52. van den Berg, I., Meuwissen, T., MacLeod, I., Goddard, M.: Predicting the effect of reference population on
799 the accuracy of within, across and multibreed genomic prediction. *J Dairy Sci* **102**(4) (2019)
- 800 53. Chen, L., Li, C., Miller, S., Schenkel, F.: Multi-population genomic prediction using a multi-task Bayesian
801 learning model. *BMC Genet* **15**, 53 (2014)
- 802 54. Esfandyari, H., Sørensen, A., Bijma, P.: A crossbred reference population can improve the response to genomic
803 selection for crossbred performance. *Genet Sel Evol* **47**(1), 76 (2015)
- 804 55. Xiang, T., Christensen, O., Legarra, A.: Technical note: Genomic evaluation for crossbred performance in a
805 single-step approach with metafounders. *J Anim Sci* **95**(4), 1472–1480 (2017)
- 806 56. Sevillano, C., Vandenplas, J., Bastiaansen, J., Bergsma, R., Calus, M.: Genomic evaluation for a three-way
807 crossbreeding system considering breed-of-origin of alleles. *Genet Sel Evol* **49**(1), 75 (2017)
- 808 57. Gautason, E., Schönherz, A., Sahana, G., Gulbrandsen, B.: Relationship of Icelandic cattle with Northern and
809 Western European cattle breeds, admixture and population structure. *Acta Agriculturae Scandinavica, Section*
810 *A-Animal Science*, 1–14 (2019)
- 811 58. Toosi, A., Fernando, R., Dekkers, J.: Genome-wide mapping of quantitative trait loci in admixed populations
812 using mixed linear model and Bayesian multiple regression analysis. *Genet Sel Evol* **50**(1), 32 (2018)
- 813 59. Vandenplas, J., Calus, M., Sevillano, C., Windig, J., Bastiaansen, J.: Assigning breed origin to alleles in
814 crossbred animals. *Genet Sel Evol* **48**(1), 61 (2016)
- 815 60. Sankararaman, S., Sridhar, S., Kimmel, G., Halperin, E.: Estimating local ancestry in admixed populations. *Am*
816 *J Hum Genet* **82**(2), 290–303 (2008)
- 817 61. Connolly, S., Fortes, M., Piper, E., Seddon, J., Kelly, M.: Determining the Number of Animals Required to
818 Accurately Determine Breed Composition Using Genomic Data (2014). 10th World Congress on Genetics
819 Applied to Livestock Production
- 820 62. Daetwyler, H., Villanueva, B., Woolliams, J.: Accuracy of predicting the genetic risk of disease using a
821 genome-wide approach. *PLoS ONE* **3**(10), 3395 (2008)
- 822 63. Wientjes, Y., Veerkamp, R., Bijma, P., Bovenhuis, H., Schrooten, C., MPL, C.: Empirical and deterministic
823 accuracies of across-population genomic prediction. *Genet Sel Evol* **47**(5) (2015)
- 824 64. Pocrnic, I., Lourenco, D., Masuda, Y., Misztal, I.: Dimensionality of genomic information and performance of
825 the algorithm for proven and young for different livestock species. *Genet Sel Evol* **48**(82) (2016)

826 Figures

Figure 1 Accuracies (horizontal axis) for high heritability trait ($h^2 = 0.4$) in QTL correlation scenario of 0.5, using different data sets or models (vertical axis). The predicted population was given on top of each plot. Letters in parenthesis stands for the significance tests.

Figure 2 Accuracies (horizontal axis) for low heritability trait ($h^2 = 0.05$) in QTL correlation scenario of 0.5, using different data sets or models (vertical axis). The predicted population was given on top of each plot. Letters in parenthesis stands for the significance tests.

Figure 3 Accuracies (horizontal axis) for high ($h^2 = 0.4$, left figure) and low ($h^2 = 0.05$, right figure) heritability trait using BOA model with correlated SNP effects and different region sizes (vertical axis), in QTL correlation scenario of 0.5. Letters in parenthesis stands for the significance tests.

Figure 4 Accuracies (horizontal axis) for low heritability trait ($h^2 = 0.05$) in QTL correlation scenario of 0.5, using different data sets or models (vertical axis), when only the QTLs are considered with region size of 1 SNP. The predicted population was given on top of each plot. Letters in parenthesis stands for the significance tests.

827 Tables

Table 1 Parents of each simulated generation

Generation/Population	HOL ¹	RED	JER	MIX
1 ²	HOL ₀ ^M × HOL ₀ ^F	RED ₀ ^M × RED ₀ ^F	JER ₀ ^M × JER ₀ ^F	JER ₀ ^M × HOL ₀ ^F
⋮	⋮	⋮	⋮	⋮
6	HOL ₅ ^M × HOL ₅ ^F	RED ₅ ^M × RED ₅ ^F	JER ₅ ^M × JER ₅ ^F	HOL ₅ ^M × MIX ₅ ^F
7	HOL ₆ ^M × HOL ₆ ^F	RED ₆ ^M × RED ₆ ^F	JER ₆ ^M × JER ₆ ^F	JER ₆ ^M × MIX ₆ ^F
8	HOL ₇ ^M × HOL ₇ ^F	RED ₇ ^M × RED ₇ ^F	JER ₇ ^M × JER ₇ ^F	RED ₇ ^M × MIX ₇ ^F
9	HOL ₈ ^M × HOL ₈ ^F	RED ₈ ^M × RED ₈ ^F	JER ₈ ^M × JER ₈ ^F	HOL ₈ ^M × MIX ₈ ^F

¹ HOL, RED, JER and MIX: Danish Holstein, Swedish Red, Danish Jersey and admixed population respectively.

² Subscripts denote the generation, and superscripts denote the sex, i.e., males (M) and females (F).

Table 2 Some descriptive statistics¹ on SNPs and QTL for each pure-breed population in the base population (Generation 0-G0)

	HOL ²	RED	JER
Number of fixed QTL for reference(alternative) allele	9(0)	5(0)	58(1)
Number of fixed SNPs for reference(alternative) allele	564(2)	385(14)	2281(286)
Number of breed specific QTL	3	4	1
Number of breed specific SNPs	261	356	50
Average MAF of segregating QTL	0.17	0.16	0.16
Average MAF of segregating SNPs	0.23	0.23	0.22

¹ average over 10 replicates

² HOL, RED and JER: Danish Holstein, Swedish Red and Danish Jersey dairy cattle, respectively

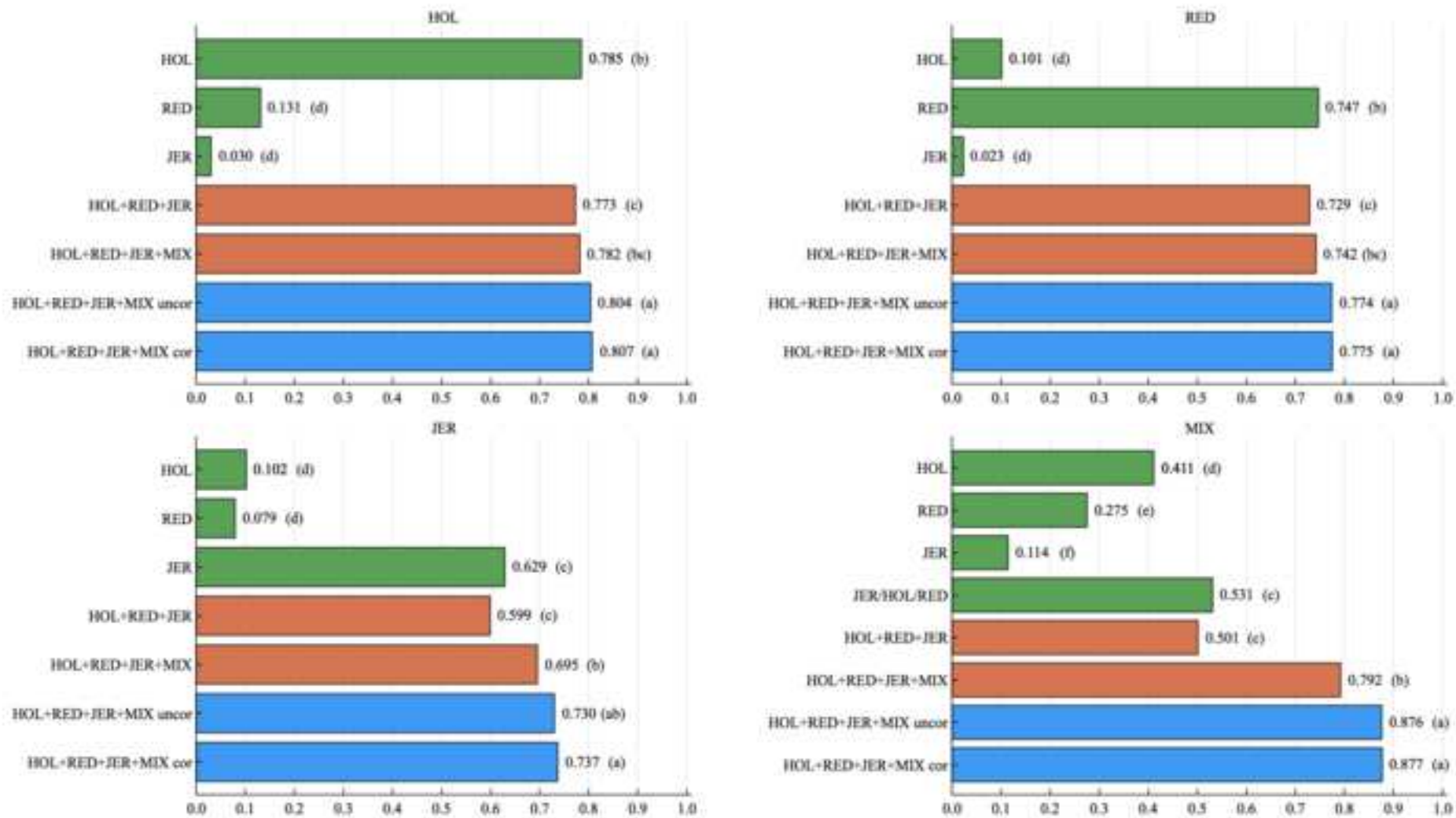
828 **Additional Files**

829 Additional File 1-Result of principle component analysis (PCA)

830 The file includes plot of the first two principle components from the PCA analysis of genomic relationship matrices.

831 Additional File 2-Results of all scenarios

832 The file includes results of all scenarios including the QTL correlation scenario of 0.5, which was partially presented
833 in the figures of main text.



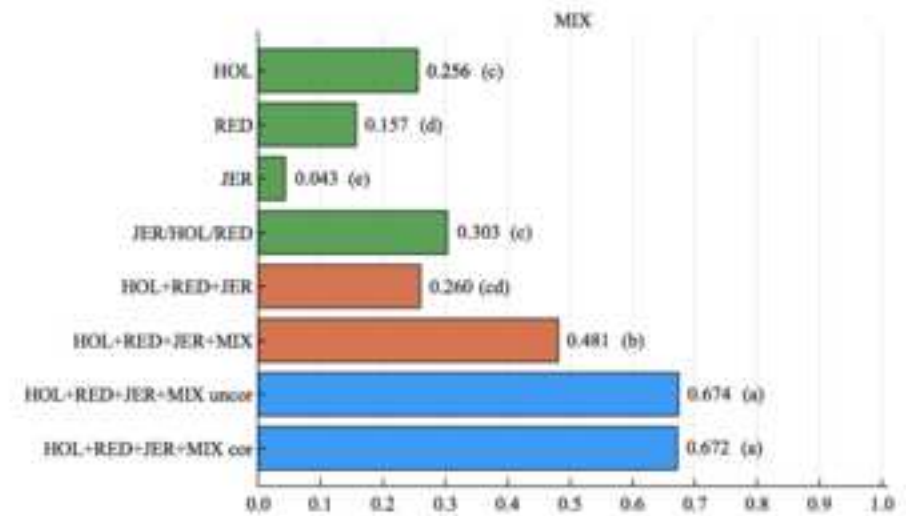
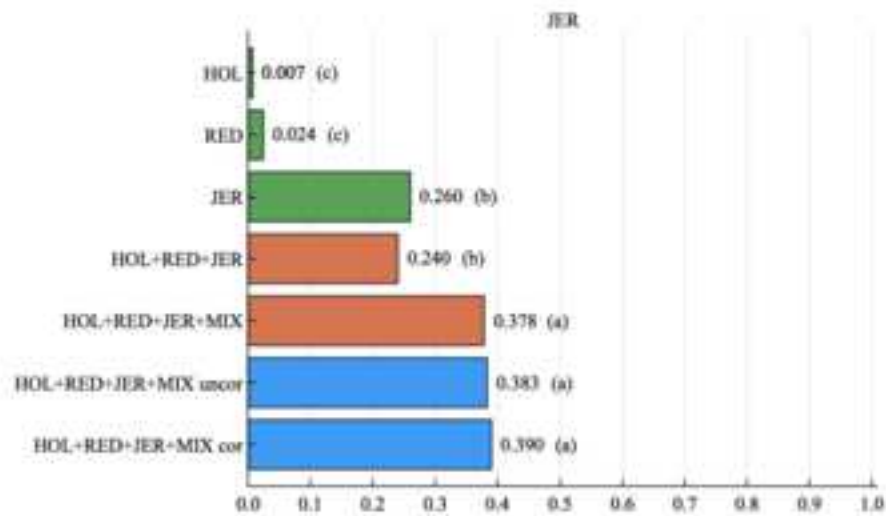
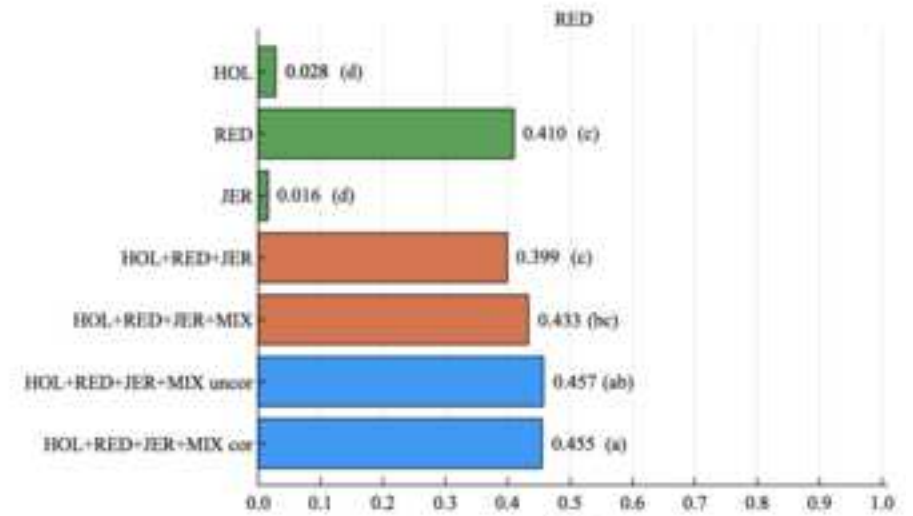
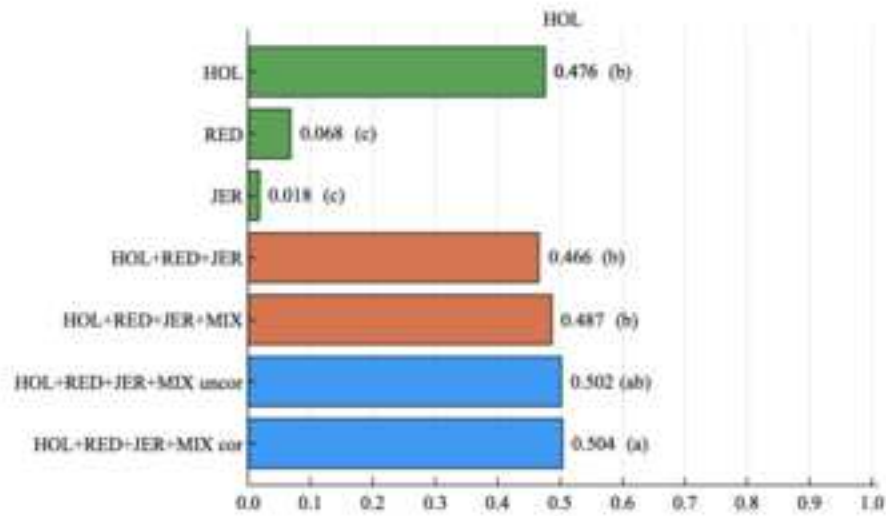


Figure 3

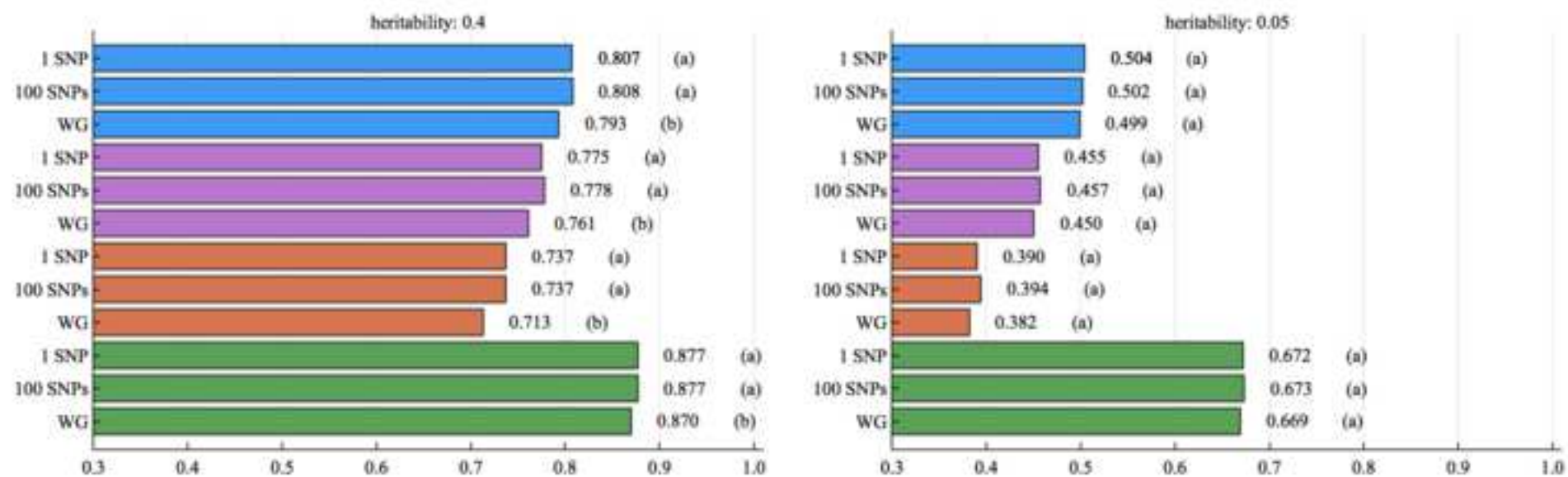
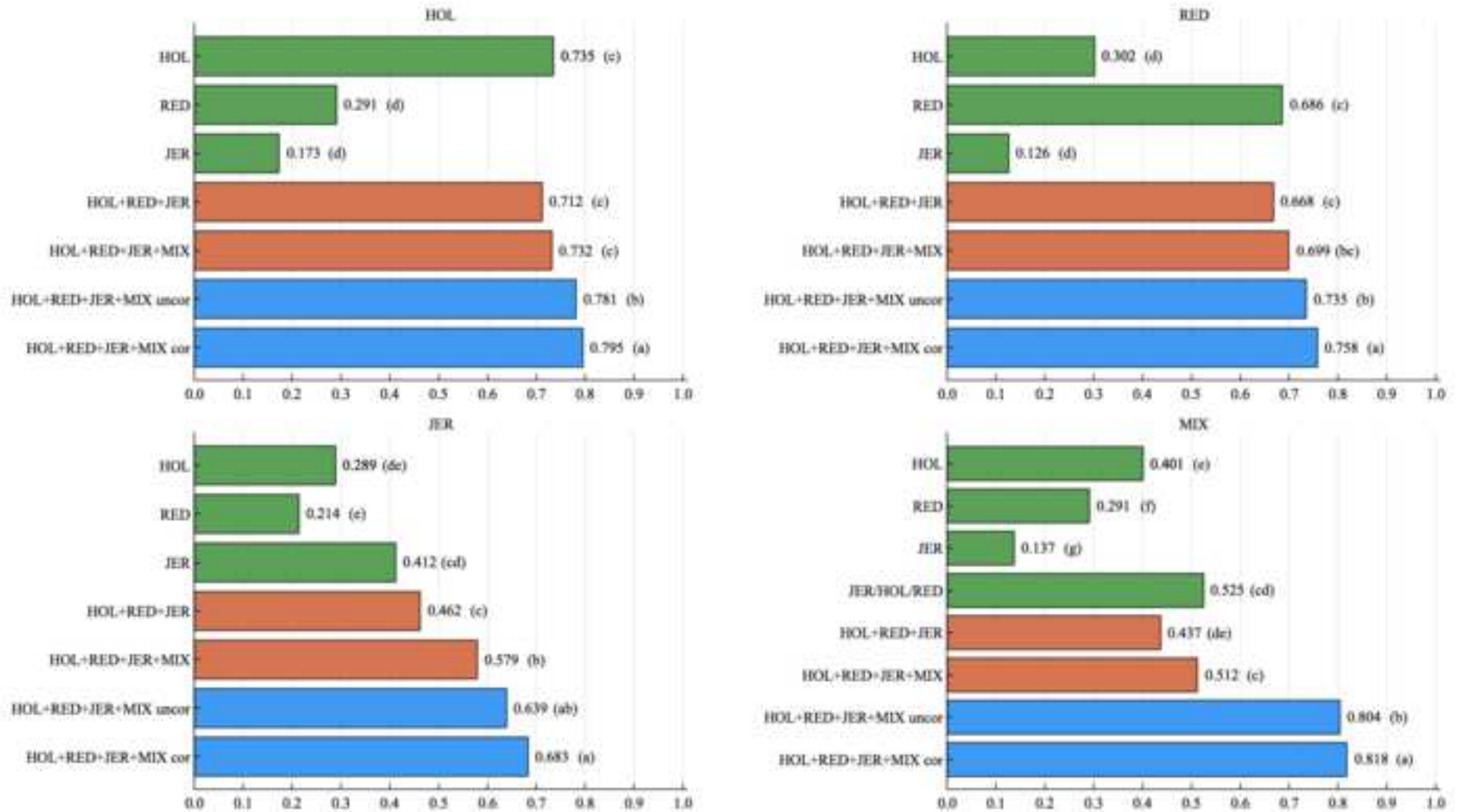
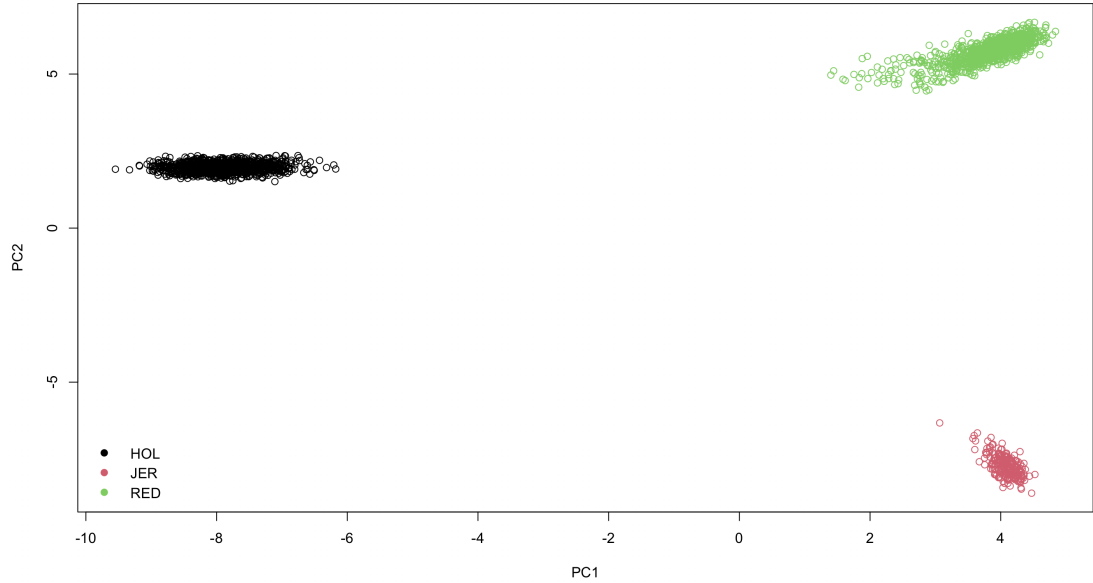


Figure 4





Supplementary Figure 1. Plot of the first two principle components from the PCA analysis of genomic relationship matrices. Genomic relationships were computed as described in [1], and analysis were carried out using R function `prcomp()` [2].

References

- [1] Wientjes, Y., Veerkamp, R., Bijma, P., Bovenhuis, H., Schrooten, C., MPL, C.: Empirical and deterministic accuracies of across-population genomic prediction. *Genet Sel Evol* **47**(5) (2015)
- [2] R Core Team: R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing, Vienna, Austria (2020). R Foundation for Statistical Computing. <https://www.R-project.org/>

Supplementary Table 1: Accuracies for purebred individuals for a trait with high heritability ($h^2 = 0.40$)

Correlation ¹	Training ²	Data/Region Size ³	HOL			RED			JER		
			1 SNP	100 SNPs	WG	1 SNP	100 SNPs	WG	1 SNP	100 SNPs	WG
1.00	Pure	HOL	b0.771 ^b	b0.780 ^a	b0.760 ^c	d0.128 ^a	d0.095 ^a	d0.080 ^b	d0.120 ^a	c0.059 ^b	c0.036 ^c
		RED	c0.170 ^a	c0.162 ^a	c0.127 ^b	c0.751 ^b	c0.757 ^a	bc0.741 ^c	d0.170 ^a	c0.123 ^b	c0.111 ^b
		JER	c0.042 ^a	d0.035 ^a	d0.040 ^a	d0.060 ^a	d0.057 ^a	d0.057 ^a	c0.643 ^b	b0.652 ^a	b0.638 ^c
	Combined	HOL+RED+JER	b0.777 ^a	b0.779 ^a	b0.757 ^b	c0.749 ^a	c0.752 ^a	c0.730 ^b	bc0.652 ^a	b0.646 ^a	b0.617 ^b
		HOL+RED+JER+MIX	a0.802 ^a	a0.803 ^a	a0.781 ^b	b0.774 ^a	b0.777 ^a	ab0.754 ^b	a0.732 ^a	a0.734 ^a	a0.701 ^b
	BOA	HOL+RED+JER+MIX no Cor	a0.796 ^b	a0.804 ^a	a0.782 ^c	ab0.776 ^b	ab0.782 ^a	a0.763 ^c	ab0.723 ^b	a0.738 ^a	a0.713 ^c
HOL+RED+JER+MIX with Cor		a0.803 ^a	a0.805 ^a	a0.783 ^b	a0.776 ^a	a0.779 ^a	a0.758 ^b	a0.734 ^a	a0.740 ^a	a0.703 ^b	
0.50	Pure	HOL	b0.785 ^b	b0.790 ^a	b0.774 ^c	d0.101 ^a	d0.085 ^{ab}	d0.073 ^b	d0.102 ^a	d0.073 ^b	e0.067 ^b
		RED	d0.131 ^{ab}	d0.129 ^a	d0.116 ^b	b0.747 ^b	b0.754 ^a	b0.736 ^c	d0.079 ^a	d0.070 ^a	e0.043 ^b
		JER	d0.030 ^a	d0.039 ^a	d0.029 ^a	d0.023 ^a	d0.025 ^a	d0.024 ^a	c0.629 ^{ab}	c0.635 ^a	cd0.624 ^b
	Combined	HOL+RED+JER	c0.773 ^a	c0.774 ^a	c0.761 ^b	c0.729 ^a	c0.732 ^a	c0.716 ^b	c0.599 ^a	c0.587 ^b	d0.590 ^{ab}
		HOL+RED+JER+MIX	bc0.782 ^a	bc0.782 ^a	bc0.771 ^b	bc0.742 ^a	bc0.747 ^a	bc0.732 ^b	b0.695 ^a	b0.690 ^a	bc0.675 ^b
	BOA	HOL+RED+JER+MIX no Cor	a0.804 ^a	a0.807 ^a	a0.792 ^b	a0.774 ^b	a0.780 ^a	a0.762 ^c	ab0.730 ^a	a0.738 ^a	ab0.712 ^b
HOL+RED+JER+MIX with Cor		a0.807 ^a	a0.808 ^a	a0.793 ^b	a0.775 ^a	a0.778 ^a	a0.761 ^b	a0.737 ^a	a0.737 ^a	a0.713 ^b	
0.25	Pure	HOL	b0.782 ^b	b0.788 ^a	b0.772 ^c	d0.053 ^a	d0.041 ^{ab}	d0.034 ^b	d0.089 ^a	d0.066 ^{ab}	d0.058 ^b
		RED	d0.080 ^a	d0.081 ^a	d0.080 ^a	b0.742 ^b	b0.750 ^a	b0.730 ^c	d0.017 ^a	d0.012 ^a	d0.003 ^a
		JER	d0.026 ^a	d0.034 ^a	d0.025 ^a	d0.017 ^a	d0.014 ^a	d0.019 ^a	bc0.632 ^a	b0.639 ^a	bc0.627 ^b
	Combined	HOL+RED+JER	c0.764 ^a	c0.766 ^a	c0.755 ^b	c0.716 ^b	c0.721 ^a	c0.705 ^c	c0.587 ^a	c0.573 ^b	c0.581 ^{ab}
		HOL+RED+JER+MIX	bc0.769 ^a	c0.768 ^{ab}	bc0.761 ^b	bc0.725 ^a	bc0.731 ^a	bc0.718 ^b	b0.689 ^a	b0.680 ^b	b0.667 ^c
	BOA	HOL+RED+JER+MIX no Cor	a0.803 ^a	a0.805 ^a	a0.790 ^b	a0.771 ^b	a0.777 ^a	a0.759 ^c	a0.735 ^a	a0.743 ^a	a0.716 ^b
HOL+RED+JER+MIX with Cor		a0.805 ^a	a0.806 ^a	a0.790 ^b	a0.772 ^a	a0.775 ^a	a0.758 ^b	a0.740 ^a	a0.741 ^a	a0.715 ^b	

¹ Correlation of simulated QTL effects. Different alphabets mean significantly different values at a Type 1 error rate of 0.05 with Bonferroni correction. Subscripts (within region size) and superscripts (within data) stand for comparisons within column and row, respectively, for each correlation scenario.

² The methods classified based on the data and model used to estimate SNP effects

³ Data: Data included in reference population. Region Size: Number of SNPs assigned the same variance.

Supplementary Table 2: Accuracies for purebred individuals for a trait with low heritability ($h^2 = 0.05$)

Correlation ¹	Training ²	Data/Region Size ³	HOL			RED			JER			
			1 SNP	100 SNPs	WG	1 SNP	100 SNPs	WG	1 SNP	100 SNPs	WG	
1.00	Pure	HOL	c0.469 ^a	c0.470 ^{ab}	a0.466 ^b	e0.019 ^a	d0.010 ^b	d0.022 ^{ab}	c-0.023 ^a	c-0.023 ^a	b-0.021 ^a	
		RED	d0.055 ^a	d0.047 ^a	b0.047 ^a	bd0.421 ^a	bc0.420 ^{ab}	bc0.417 ^b	c0.044 ^a	c0.043 ^a	b0.045 ^a	
		JER	d-0.008 ^a	d-0.002 ^a	b0.002 ^a	e0.036 ^a	d0.035 ^a	d0.040 ^a	ab0.319 ^a	ab0.319 ^a	a0.321 ^a	
	Combined	HOL+RED+JER	bc0.467 ^a	c0.472 ^a	a0.463 ^a	cd0.415 ^a	c0.411 ^a	c0.408 ^a	b0.291 ^a	b0.294 ^a	a0.289 ^a	
		HOL+RED+JER+MIX	ab0.501 ^a	b0.503 ^{ab}	a0.494 ^b	ab0.453 ^a	ab0.450 ^{ab}	ab0.445 ^b	ab0.389 ^a	a0.385 ^a	a0.381 ^a	
	BOA	HOL+RED+JER+MIX no Cor	abc0.500 ^a	abc0.501 ^a	a0.497 ^a	ac0.451 ^a	a0.451 ^{ab}	a0.447 ^b	ab0.387 ^a	ab0.382 ^a	a0.383 ^a	
		HOL+RED+JER+MIX with Cor	a0.503 ^a	a0.506 ^{ab}	a0.495 ^b	ab0.454 ^a	ab0.450 ^{ab}	ab0.444 ^b	a0.392 ^a	a0.389 ^a	a0.382 ^a	
	0.50	Pure	HOL	b0.476 ^a	a0.478 ^a	a0.474 ^a	d0.028 ^a	c0.024 ^a	c0.030 ^a	c0.007 ^a	c0.017 ^a	c0.010 ^a
			RED	c0.068 ^a	b0.080 ^a	b0.071 ^a	c0.410 ^a	b0.412 ^a	b0.410 ^a	c0.024 ^a	c0.030 ^a	c0.028 ^a
JER			c0.018 ^a	b0.015 ^a	b0.012 ^a	d0.016 ^{ab}	c0.013 ^b	c0.026 ^a	b0.260 ^a	b0.260 ^a	b0.264 ^a	
Combined		HOL+RED+JER	b0.466 ^a	a0.467 ^a	a0.464 ^a	c0.399 ^a	b0.393 ^a	b0.397 ^a	b0.240 ^a	b0.238 ^a	b0.243 ^a	
		HOL+RED+JER+MIX	b0.487 ^a	a0.483 ^a	a0.482 ^a	bc0.433 ^a	b0.432 ^a	ab0.432 ^a	a0.378 ^a	a0.374 ^a	a0.372 ^a	
BOA		HOL+RED+JER+MIX no Cor	ab0.502 ^a	a0.498 ^a	a0.499 ^a	ab0.457 ^a	a0.459 ^a	a0.455 ^a	a0.383 ^a	a0.388 ^a	a0.384 ^a	
		HOL+RED+JER+MIX with Cor	a0.504 ^a	a0.502 ^a	a0.499 ^a	a0.455 ^a	a0.457 ^a	a0.450 ^a	a0.390 ^a	a0.394 ^a	a0.382 ^a	
0.25		Pure	HOL	ab0.473 ^a	a0.474 ^a	ab0.471 ^a	c0.012 ^a	d0.007 ^a	e0.014 ^a	c0.005 ^a	c0.016 ^a	c0.007 ^a
			RED	c0.050 ^a	b0.059 ^a	c0.054 ^a	b0.407 ^a	b0.409 ^a	c0.407 ^a	c-0.006 ^a	c-0.006 ^a	c0.000 ^a
	JER		c0.017 ^a	b0.014 ^a	c0.012 ^a	c0.024 ^b	d0.017 ^b	e0.032 ^a	b0.260 ^a	b0.260 ^a	b0.265 ^a	
	Combined	HOL+RED+JER	b0.459 ^a	a0.458 ^a	b0.458 ^a	b0.392 ^a	c0.385 ^{ab}	d0.391 ^a	b0.228 ^a	b0.226 ^a	b0.233 ^a	
		HOL+RED+JER+MIX	b0.476 ^a	a0.472 ^a	b0.473 ^a	b0.423 ^a	bc0.417 ^a	bcd0.424 ^a	a0.372 ^a	a0.364 ^a	a0.368 ^a	
	BOA	HOL+RED+JER+MIX no Cor	a0.499 ^a	a0.494 ^a	a0.495 ^a	a0.454 ^b	a0.458 ^a	ab0.453 ^{ab}	a0.388 ^a	a0.392 ^a	a0.389 ^a	
		HOL+RED+JER+MIX with Cor	a0.501 ^a	a0.497 ^a	ab0.495 ^a	a0.454 ^{ab}	a0.458 ^a	a0.451 ^b	a0.393 ^a	a0.397 ^a	a0.384 ^a	

¹ Correlation of simulated QTL effects. Different alphabets mean significantly different values at a Type 1 error rate of 0.05 with Bonferroni correction. Subscripts (within region size) and superscripts (within data) stand for comparisons within column and row, respectively, for each correlation scenario.

² The methods classified based on the data and model used to estimate SNP effects

³ Data: Data included in reference population. Region Size: Number of SNPs assigned the same variance.

Supplementary Table 3: Accuracies for admixed individuals for a trait with high heritability ($h^2 = 0.40$)

Correlation ¹	Training ²	Data/Region Size ³	MIX		
			1 SNP	100 SNPs	WG
1.00	Pure	HOL	_e 0.578 ^a	_e 0.576 ^a	_e 0.562 ^b
		RED	_f 0.395 ^a	_f 0.395 ^a	_f 0.372 ^b
		JER	_g 0.178 ^a	_g 0.179 ^a	_g 0.175 ^a
		HOL/RED/JER	_d 0.687 ^b	_d 0.704 ^a	_d 0.685 ^b
	Combined	HOL+RED+JER	_c 0.726 ^a	_c 0.725 ^a	_c 0.704 ^b
		HOL+RED+JER+MIX	_a 0.788 ^a	_a 0.792 ^a	_{ab} 0.770 ^b
	BOA	HOL+RED+JER+MIX no Cor	_b 0.772 ^b	_b 0.782 ^a	_b 0.762 ^c
		HOL+RED+JER+MIX with Cor	_a 0.794 ^b	_a 0.799 ^a	_a 0.776 ^c
0.50	Pure	HOL	_d 0.411 ^a	_d 0.415 ^a	_d 0.410 ^a
		RED	_e 0.275 ^a	_e 0.274 ^a	_e 0.268 ^a
		JER	_f 0.114 ^a	_f 0.117 ^a	_f 0.114 ^a
		HOL/RED/JER	_c 0.531 ^a	_c 0.533 ^{ab}	_c 0.521 ^b
	Combined	HOL+RED+JER	_c 0.501 ^a	_c 0.498 ^{ab}	_c 0.493 ^b
		HOL+RED+JER+MIX	_b 0.792 ^a	_b 0.783 ^a	_b 0.774 ^b
	BOA	HOL+RED+JER+MIX no Cor	_a 0.876 ^a	_a 0.877 ^a	_a 0.870 ^b
		HOL+RED+JER+MIX with Cor	_a 0.877 ^a	_a 0.877 ^a	_a 0.870 ^b
0.25	Pure	HOL	_d 0.358 ^a	_d 0.363 ^a	_d 0.360 ^a
		RED	_e 0.231 ^a	_e 0.231 ^a	_e 0.231 ^a
		JER	_e 0.106 ^a	_e 0.109 ^a	_e 0.106 ^a
		HOL/RED/JER	_c 0.482 ^a	_c 0.482 ^{ab}	_c 0.469 ^b
	Combined	HOL+RED+JER	_c 0.434 ^a	_c 0.433 ^{ab}	_c 0.430 ^b
		HOL+RED+JER+MIX	_b 0.793 ^{ab}	_b 0.788 ^a	_b 0.777 ^b
	BOA	HOL+RED+JER+MIX no Cor	_a 0.897 ^a	_a 0.898 ^a	_a 0.892 ^b
		HOL+RED+JER+MIX with Cor	_a 0.897 ^a	_a 0.897 ^a	_a 0.892 ^b

¹ Correlation of simulated QTL effects. Different alphabets mean significantly different values at a Type 1 error rate of 0.05 with Bonferroni correction. Subscripts (within region size) and superscripts (within data) stand for comparisons within column and row, respectively, for each correlation scenario.

² The methods classified based on the data and model used to estimate SNP effects

³ Data: Data included in reference population. Region Size: Number of SNPs assigned the same variance.

Supplementary Table 4: Accuracies for admixed individuals for a trait with low heritability ($h^2 = 0.05$)

Correlation ¹	Training ²	Data/Region Size ³	MIX		
			1 SNP	100 SNPs	WG
1.00	Pure	HOL	c0.363 ^a	d0.364 ^a	d0.362 ^a
		RED	d0.229 ^a	e0.225 ^{ab}	e0.221 ^b
		JER	e0.075 ^a	f0.077 ^a	f0.078 ^a
		HOL/RED/JER	b0.433 ^a	bc0.436 ^a	bc0.429 ^a
	Combined	HOL+RED+JER	b0.436 ^a	c0.435 ^a	c0.429 ^a
		HOL+RED+JER+MIX	a0.491 ^a	ab0.488 ^a	ab0.484 ^a
	BOA	HOL+RED+JER+MIX no Cor	bc0.430 ^{ab}	cd0.434 ^a	cd0.424 ^b
		HOL+RED+JER+MIX with Cor	a0.500 ^a	a0.499 ^a	a0.493 ^a
0.50	Pure	HOL	c0.256 ^a	c0.256 ^a	c0.255 ^a
		RED	d0.157 ^a	d0.162 ^a	de0.155 ^a
		JER	e0.043 ^a	e0.043 ^a	e0.043 ^a
		HOL/RED/JER	c0.303 ^{ab}	c0.308 ^a	c0.299 ^b
	Combined	HOL+RED+JER	cd0.260 ^a	c0.260 ^a	cd0.258 ^a
		HOL+RED+JER+MIX	b0.481 ^a	b0.471 ^b	b0.477 ^{ab}
	BOA	HOL+RED+JER+MIX no Cor	a0.674 ^{ab}	a0.676 ^a	a0.672 ^b
		HOL+RED+JER+MIX with Cor	a0.672 ^a	a0.673 ^a	a0.669 ^a
0.25	Pure	HOL	cd0.226 ^a	cd0.227 ^a	cd0.226 ^a
		RED	de0.136 ^a	de0.138 ^a	de0.134 ^a
		JER	e0.040 ^a	e0.041 ^a	e0.041 ^a
		HOL/RED/JER	c0.271 ^{ab}	c0.277 ^a	c0.266 ^b
	Combined	HOL+RED+JER	cd0.222 ^a	cd0.222 ^a	cd0.221 ^a
		HOL+RED+JER+MIX	b0.494 ^a	b0.485 ^a	b0.490 ^a
	BOA	HOL+RED+JER+MIX no Cor	a0.729 ^a	a0.730 ^{ab}	a0.727 ^b
		HOL+RED+JER+MIX with Cor	a0.728 ^a	a0.728 ^a	a0.725 ^a

¹ Correlation of simulated QTL effects. Different alphabets mean significantly different values at a Type 1 error rate of 0.05 with Bonferroni correction. Subscripts (within region size) and superscripts (within data) stand for comparisons within column and row, respectively, for each correlation scenario.

² The methods classified based on the data and model used to estimate SNP effects

³ Data: Data included in reference population. Region Size: Number of SNPs assigned the same variance.

Supplementary Table 5: Accuracy for low heritability ($h^2 = 0.05$) trait, using only 250 QTL and region size of 1 SNP

Correlation ¹	Training ²	Data/Region Size ³	HOL	RED	JER	MIX
			1 SNP	1 SNP	1 SNP	1 SNP
0.50	Pure	HOL	_c 0.735	_d 0.302	_{de} 0.289	_e 0.401
		RED	_d 0.291	_c 0.686	_e 0.214	_f 0.291
		JER	_d 0.173	_d 0.126	_{cd} 0.412	_g 0.137
		HOL/RED/JER				_{cd} 0.525
	Combined	HOL+RED+JER	_c 0.712	_c 0.668	_c 0.462	_{de} 0.437
		HOL+RED+JER+MIX	_c 0.732	_{bc} 0.699	_b 0.579	_c 0.512
	BOA	HOL+RED+JER+MIX no Cor	_b 0.781	_b 0.735	_{ab} 0.639	_b 0.804
		HOL+RED+JER+MIX with Cor	_a 0.795	_a 0.758	_a 0.683	_a 0.818

¹ Correlation of simulated QTL effects. Different alphabets mean significantly different values at a Type 1 error rate of 0.05 with Bonferroni correction. Subscripts (region size of 1 SNP) and superscripts (within data) stand for comparisons within column and row, respectively, for each correlation scenario.

² The methods classified based on the data and model used to estimate SNP effects

³ Data: Data included in reference population.