



Genomic prediction using individual-level data and summary statistics from multiple populations

Vandenplas, J., Calus, M. P. L., & Gorjanc, G.

This is a "Post-Print" accepted manuscript, which has been published in "Genetics"

This version is distributed under a non-commercial no derivatives Creative Commons



([CC-BY-NC-ND](https://creativecommons.org/licenses/by-nc-nd/4.0/)) user license, which permits use, distribution, and reproduction in any medium, provided the original work is properly cited and not used for commercial purposes. Further, the restriction applies that if you remix, transform, or build upon the material, you may not distribute the modified material.

Please cite this publication as follows:

Vandenplas, J., Calus, M. P. L., & Gorjanc, G. (2018). Genomic prediction using individual-level data and summary statistics from multiple populations. *Genetics*, 210(1), 53-69. DOI: 10.1534/genetics.118.301109

You can download the published version at:

<https://doi.org/10.1534/genetics.118.301109>

1 **Genomic prediction using individual-level data and summary statistics from multiple**
2 **populations**

3 Jeremie Vandenplas*, Mario P.L. Calus*, Gregor Gorjanc†

4

5 * Wageningen University & Research, Animal Breeding and Genomics, 6700 AH

6 Wageningen, The Netherlands

7 † The Roslin Institute and Royal (Dick) School of Veterinary Studies, University of

8 Edinburgh, Easter Bush Research Centre, Midlothian EH25 9RG, UK

9

Multi-population genomic prediction

10 Running title: Multi-population genomic prediction

11

12 Key words: meta-analysis, quantitative trait, statistical method

13

14 Author information:

15 Jeremie Vandenplas

16 Wageningen University & Research

17 Animal Breeding and Genomics

18 P.O. box 338, 6700 AH Wageningen, the Netherlands

19 E-mail: [jeremie.vandenplas @wur.nl](mailto:jeremie.vandenplas@wur.nl)

20 Phone: +31 06 83642304

21

22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42

ABSTRACT

This study presents a method for genomic prediction that uses individual-level data and summary statistics from multiple populations. Genome-wide markers are nowadays widely used to predict complex traits, and genomic prediction using multi-population data is an appealing approach to achieve higher prediction accuracies. However, sharing of individual-level data across populations is not always possible. We present a method that enables integration of summary statistics from separate analyses with the available individual-level data. The data can either consist of individuals with single or multiple (weighted) phenotype records per individual. We developed a method based on a hypothetical joint analysis model and absorption of population specific information. We show that population specific information is fully captured by estimated allele substitution effects and the accuracy of those estimates, i.e. the summary statistics. The method gives identical result as the joint analysis of all individual-level data when complete summary statistics are available. We provide a series of easy-to-use approximations that can be used when complete summary statistics are not available or impractical to share. Simulations show that approximations enables integration of different sources of information across a wide range of settings yielding accurate predictions. The method can be readily extended to multiple-traits. In summary, the developed method enables integration of genome-wide data in the individual-level or summary statistics form from multiple populations to obtain more accurate estimates of allele substitution effects and genomic predictions.

43

INTRODUCTION

44 Genome-wide markers are nowadays widely used in animal and plant breeding to
45 predict complex traits. This prediction is based on a linear model that partitions for each
46 individual the observed complex phenotype value into systematic effects, comprising at least a
47 population mean, an individual genetic value and an environmental deviation (Fisher, 1918).
48 With genome-wide markers, individual genetic values can be computed from allele substitution
49 effects estimated from individual-level phenotype and genotype data (Meuwissen et al., 2001).
50 Subsequently, genetic values can be also computed for individuals of interest that are
51 genotyped, but not phenotyped. This process is commonly called genomic prediction. In animal
52 and plant breeding, genetic values are used to identify genetically superior individuals and use
53 them as parents of the next generation to improve complex traits like milk yield (Meuwissen et
54 al., 2001; VanRaden, 2008) or grain yield (Schulthess et al., 2016). In human genetics, genetic
55 values can be used to predict individual genetic risk for complex diseases to inform preventive
56 and personalized medicine (Campos et al., 2010; Wray et al., 2013; Pasaniuc and Price, 2017).

57 Accuracy of estimated allele substitution effects and of resulting genetic values for
58 complex traits are foremost a function of the number of individuals with available phenotypes
59 and genotypes (Daetwyler et al., 2008). To maximize the prediction accuracy, use of all
60 available data is recommended (Henderson, 1984; Wray et al., 2013; Vilhjálmsson et al., 2015).
61 In some small populations, collecting large amounts of data is not possible, and a joint analysis
62 across multiple populations is needed to achieve high accuracy (Hozé et al., 2014; Wientjes et
63 al., 2016). However, such joint analysis is often impossible, because of logistic or privacy
64 considerations (Powell and Norman, 1998; Maier et al., 2018). Therefore, several methods were
65 proposed to enable analysis of data from multiple populations when individual-level data is not
66 available (Pasaniuc and Price, 2017; Liu and Goddard, 2018; Maier et al., 2018). These
67 methods, often called meta-analyses (Pasaniuc and Price, 2017), approximate a joint analysis

68 by first obtaining summary statistics from separate analyses of individual-level data for each
69 population and then combine these summary statistics to estimate genetic values. In human
70 genetics, summary statistics usually consist of publically available allele substitution effects,
71 i.e., genome-wide associations, together with their standard errors, estimated independently for
72 each marker (Yang et al., 2012; Vilhjálmsson et al., 2015; Maier et al., 2018). In livestock,
73 summary statistics more likely consist of allele substitution effects estimated jointly for all
74 markers, together with prediction error (co)variances (Liu and Goddard, 2018). While these
75 methods may increase prediction accuracy in comparison to separate analyses, a loss in
76 prediction accuracy is expected relative to an analysis using all individual-level data due to
77 approximations (Maier et al., 2018). Further, these methods are based on some assumptions that
78 make them difficult to apply outside their context of development. For example, Maier et al.
79 (2018) implicitly assumed that only a single phenotype record per trait was associated with an
80 individual. While this is usually the case in human genetics, it is not in breeding populations
81 where individuals may have repeated phenotype records for the same trait, e.g., repeated
82 longitudinal production or reproduction records in livestock or replicated field trials in crops,
83 or when phenotype records are measured on a group of individuals and linked to a genotyped
84 relative, e.g., progeny tested bulls for dairy production. Also, these developed methods do not
85 allow combining individual-level data from some and summary statistics from other
86 populations in one analysis (Liu and Goddard, 2018; Maier et al., 2018).

87 The objective of this study was to develop a method that jointly analyses individual-
88 level data and summary statistics from multiple populations with no or limited amount of
89 approximation. The method assumes that individual-level data is composed of marker
90 genotypes and phenotype records that potentially have a variable number of replicates per
91 individual. Further, summary statistics are assumed to be composed of estimated allele
92 substitution effects with an associated measure of accuracy. Different measures of accuracy can

Multi-population genomic prediction

93 be used, which controls the amount of approximation. The developed method is validated with
94 simulated data. The results show that the method enables accurate integration of different
95 sources of information across a wide range of settings.

96

97

MATERIAL AND METHODS

98

99

100

101

102

The first part of this section describes the theory of (1) separate and joint analyses of two individual-level datasets, (2) an exact integration of estimated allele substitution effects from one population into the analysis of another, (3) approximate integrations, and (4) generalization for multiple populations. The second part describes simulations used for validation of the developed method.

103

Theory

104

105

106

107

108

109

110

111

Assume we have two populations with independent individual-level datasets of phenotyped and genotyped individuals. The two populations and their corresponding datasets are hereafter referred to as 1 and 2. Further assume that both datasets contain the same markers. From this data we want to obtain accurate estimates of allele substitution effects and genetic values for complex traits. We can achieve this by a joint analysis of the two datasets. When one of the datasets is not available, we can achieve this by integrating the results of a separate analysis of the unavailable data into the separate analysis of the available dataset. We show how to perform this integration exactly or approximately.

112

Separate and joint analyses

113

114

A standard marker model, using random regression on marker genotypes, for the separate analysis of dataset i ($i = 1, 2$) is:

115

$$\mathbf{y}_i = \mathbf{X}_i \boldsymbol{\beta}_i^* + \mathbf{Z}_i \mathbf{W}_i \boldsymbol{\alpha}_i^* + \mathbf{e}_i^*, \quad (1)$$

116

117

118

where \mathbf{y}_i is a $n_{obs,i} \times 1$ vector of phenotypes, $\boldsymbol{\beta}_i^*$ is a $n_{f,i} \times 1$ vector of fixed effects that are linked to \mathbf{y}_i by a $n_{obs,i} \times n_{f,i}$ incidence matrix \mathbf{X}_i , $\boldsymbol{\alpha}_i^*$ is a $n_{mar} \times 1$ vector of allele substitution effects that are linked to \mathbf{y}_i by a $n_{obs,i} \times n_{ind,i}$ incidence matrix \mathbf{Z}_i and a $n_{ind,i} \times$

Multi-population genomic prediction

119 n_{mar} matrix of genotypes \mathbf{W}_i , and \mathbf{e}_i^* is the vector $n_{obs,i} \times 1$ of residuals. In this work we
 120 consider single-nucleotide polymorphism markers, which we code in \mathbf{W}_i as 0 for homozygous
 121 aa, 1 for heterozygous aA or Aa, and 2 for homozygous AA. Other genotype coding and
 122 centering, that is of the form $(\mathbf{W}_i - \mathbf{1}\mathbf{v}_i')$ with $\mathbf{1}$ being a $n_{ind,i} \times 1$ vector of ones and \mathbf{v}_i being a
 123 $n_{mar} \times 1$ vector, can be used with no difference in obtained estimates of allele substitution
 124 effects (Strandén and Christensen, 2011). We assume a prior multivariate normal (MVN)
 125 distribution for allele substitution effects for the separate analysis of the dataset i , $\boldsymbol{\alpha}_i^*$, with mean
 126 zero and covariance $\mathbf{B}_i \sigma_{\alpha_i}^2$, $\boldsymbol{\alpha}_i^* \sim MVN(\mathbf{0}, \mathbf{B}_i \sigma_{\alpha_i}^2)$, where \mathbf{B}_i is a $n_{mar} \times n_{mar}$ diagonal matrix
 127 (e.g., an identity matrix \mathbf{I}), and $\sigma_{\alpha_i}^2$ is the variance of allele substitution effects. We also assume
 128 that residuals are multivariate normally distributed with mean zero and covariance $\mathbf{R}_i \sigma_e^2$,
 129 $\mathbf{e}_i^* \sim MVN(\mathbf{0}, \mathbf{R}_i \sigma_e^2)$, where \mathbf{R}_i is a $n_{obs,i} \times n_{obs,i}$ diagonal matrix (e.g., an identity matrix \mathbf{I}),
 130 and σ_e^2 is the residual variance. For simplicity and without loss of generality, it is assumed in
 131 the following that residual variances are the same for all separate and joint analyses. Variance
 132 components $\sigma_{\alpha_i}^2$ and σ_e^2 are assumed known, as they will have been estimated from the data
 133 previously. This marker model is the ridge regression model (Hoerl and Kennard, 1976;
 134 Whittaker et al., 2000; Meuwissen et al., 2001; de los Campos et al., 2012) with optional
 135 different weights in \mathbf{B}_i (to differentially shrink different loci) and \mathbf{R}_i (to account for
 136 heterogeneous residual variance due to variable number of repeated phenotype records per
 137 individual).

138 Separate estimates of allele substitution effects $\widehat{\boldsymbol{\alpha}}_i^*$ are obtained by solving the following
 139 system of equations:

$$140 \quad \begin{bmatrix} \mathbf{X}_i' \mathbf{R}_i^{-1} \sigma_e^{-2} \mathbf{X}_i & \mathbf{X}_i' \mathbf{R}_i^{-1} \sigma_e^{-2} \mathbf{Z}_i \mathbf{W}_i \\ \mathbf{W}_i' \mathbf{Z}_i' \mathbf{R}_i^{-1} \sigma_e^{-2} \mathbf{X}_i & \mathbf{W}_i' \mathbf{Z}_i' \mathbf{R}_i^{-1} \sigma_e^{-2} \mathbf{Z}_i \mathbf{W}_i + \mathbf{B}_i^{-1} \sigma_{\alpha_i}^{-2} \end{bmatrix} \begin{bmatrix} \widehat{\boldsymbol{\beta}}_i^* \\ \widehat{\boldsymbol{\alpha}}_i^* \end{bmatrix} = \begin{bmatrix} \mathbf{X}_i' \mathbf{R}_i^{-1} \sigma_e^{-2} \mathbf{y}_i \\ \mathbf{W}_i' \mathbf{Z}_i' \mathbf{R}_i^{-1} \sigma_e^{-2} \mathbf{y}_i \end{bmatrix}. \quad (2)$$

141 Separate estimates of genetic values for individuals in a dataset i ($i = 1, 2$) are

142 obtained by $\widehat{\mathbf{g}}_i^* = \mathbf{W}_i \widehat{\boldsymbol{\alpha}}_i^*$.

143 A marker model for the joint analysis of two datasets 1 and 2 is:

$$144 \begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{X}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{X}_2 \end{bmatrix} \begin{bmatrix} \boldsymbol{\beta}_1 \\ \boldsymbol{\beta}_2 \end{bmatrix} + \begin{bmatrix} \mathbf{Z}_1 & \mathbf{W}_1 \\ \mathbf{Z}_2 & \mathbf{W}_2 \end{bmatrix} \boldsymbol{\alpha} + \begin{bmatrix} \mathbf{e}_1 \\ \mathbf{e}_2 \end{bmatrix}, \quad (3)$$

145 where phenotypes from the two populations are modelled with populations specific fixed effects

146 $(\boldsymbol{\beta}_1, \boldsymbol{\beta}_2)$, but a joint set of allele substitution effects $(\boldsymbol{\alpha})$. We assume a multivariate normal

147 prior distribution for allele substitution effects with mean zero and covariance $\mathbf{B}_J \sigma_{\alpha_j}^2$,

148 $\boldsymbol{\alpha} \sim MVN(\mathbf{0}, \mathbf{B}_J \sigma_{\alpha_j}^2)$, where \mathbf{B}_J is a $n_{mar} \times n_{mar}$ diagonal matrix, and $\sigma_{\alpha_j}^2$ is the variance of

149 allele substitution effects in the joint analysis. We also assume that residuals are multivariate

150 normally distributed, specifically $\begin{bmatrix} \mathbf{e}_1 \\ \mathbf{e}_2 \end{bmatrix} \sim MVN\left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \mathbf{R}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{R}_2 \end{bmatrix} \sigma_e^2\right)$ where \mathbf{R}_i is a $n_{obs,i} \times$

151 $n_{obs,i}$ diagonal matrix.

152 Joint estimates of allele substitution effects $\widehat{\boldsymbol{\alpha}}$ are obtained by solving the following

153 system of equations:

$$154 \begin{bmatrix} \mathbf{X}'_1 \mathbf{R}_1^{-1} \sigma_e^{-2} \mathbf{X}_1 & \mathbf{0} & \mathbf{X}'_1 \mathbf{R}_1^{-1} \sigma_e^{-2} \mathbf{Z}_1 \mathbf{W}_1 \\ \mathbf{0} & \mathbf{X}'_2 \mathbf{R}_2^{-1} \sigma_e^{-2} \mathbf{X}_2 & \mathbf{X}'_2 \mathbf{R}_2^{-1} \sigma_e^{-2} \mathbf{Z}_2 \mathbf{W}_2 \\ \mathbf{W}'_1 \mathbf{Z}'_1 \mathbf{R}_1^{-1} \sigma_e^{-2} \mathbf{X}_1 & \mathbf{W}'_2 \mathbf{Z}'_2 \mathbf{R}_2^{-1} \sigma_e^{-2} \mathbf{X}_2 & \mathbf{W}'_1 \mathbf{Z}'_1 \mathbf{R}_1^{-1} \sigma_e^{-2} \mathbf{Z}_1 \mathbf{W}_1 + \mathbf{W}'_2 \mathbf{Z}'_2 \mathbf{R}_2^{-1} \sigma_e^{-2} \mathbf{Z}_2 \mathbf{W}_2 + \mathbf{B}_J^{-1} \sigma_{\alpha_j}^{-2} \end{bmatrix} \begin{bmatrix} \widehat{\boldsymbol{\beta}}_1 \\ \widehat{\boldsymbol{\beta}}_2 \\ \widehat{\boldsymbol{\alpha}} \end{bmatrix} =$$

$$155 \begin{bmatrix} \mathbf{X}'_1 \mathbf{R}_1^{-1} \sigma_e^{-2} \mathbf{y}_1 \\ \mathbf{X}'_2 \mathbf{R}_2^{-1} \sigma_e^{-2} \mathbf{y}_2 \\ \mathbf{W}'_1 \mathbf{Z}'_1 \mathbf{R}_1^{-1} \sigma_e^{-2} \mathbf{y}_1 + \mathbf{W}'_2 \mathbf{Z}'_2 \mathbf{R}_2^{-1} \sigma_e^{-2} \mathbf{y}_2 \end{bmatrix} \quad (4).$$

156 Joint estimates of genetic values for individuals in a dataset i ($i = 1, 2$) are obtained by

157 $\widehat{\mathbf{g}}_i = \mathbf{W}_i \widehat{\boldsymbol{\alpha}}$.

158 *Exact integration*

159 The integration of estimates of allele substitution effects from one dataset into the
 160 analysis of another can be performed by means of absorbing corresponding equations in the
 161 joint system of equations. We choose to integrate estimates from the dataset 1 into the analysis
 162 of dataset 2. Derivations in Appendix A1 lead to the following system of equations that
 163 performs such integration and gives equivalent estimates of allele substitution effects to the
 164 joint analysis (4):

$$\begin{aligned}
 & \begin{bmatrix} \mathbf{X}'_2 \mathbf{R}_2^{-1} \sigma_e^{-2} \mathbf{X}_2 & \mathbf{X}'_2 \mathbf{R}_2^{-1} \sigma_e^{-2} \mathbf{Z}_2 \mathbf{W}_2 \\ \mathbf{W}'_2 \mathbf{Z}'_2 \mathbf{R}_2^{-1} \sigma_e^{-2} \mathbf{X}_2 & \left(PEC(\widehat{\boldsymbol{\alpha}}_1^*)\right)^{-1} + \mathbf{W}'_2 \mathbf{Z}'_2 \mathbf{R}_2^{-1} \sigma_e^{-2} \mathbf{Z}_2 \mathbf{W}_2 - \mathbf{B}_1^{-1} \sigma_{\alpha_1}^{-2} + \mathbf{B}_J^{-1} \sigma_{\alpha_J}^{-2} \end{bmatrix} \begin{bmatrix} \widehat{\boldsymbol{\beta}}_2 \\ \widehat{\boldsymbol{\alpha}} \end{bmatrix} = \\
 & \begin{bmatrix} \mathbf{X}'_2 \mathbf{R}_2^{-1} \sigma_e^{-2} \mathbf{y}_2 \\ \left(PEC(\widehat{\boldsymbol{\alpha}}_1^*)\right)^{-1} \widehat{\boldsymbol{\alpha}}_1^* + \mathbf{W}'_2 \mathbf{Z}'_2 \mathbf{R}_2^{-1} \sigma_e^{-2} \mathbf{y}_2 \end{bmatrix}, \tag{5}
 \end{aligned}$$

167 where $\widehat{\boldsymbol{\alpha}}_1^*$ are estimates of allele substitution effects from the separate analysis of dataset 1 using
 168 (2), and $\left(PEC(\widehat{\boldsymbol{\alpha}}_1^*)\right)^{-1}$ is the inverse of the corresponding prediction error covariance (PEC)
 169 matrix. The latter can be obtained as $\left(PEC(\widehat{\boldsymbol{\alpha}}_1^*)\right)^{-1} = \mathbf{W}'_1 \mathbf{Z}'_1 \mathbf{M}_1 \sigma_e^{-2} \mathbf{Z}_1 \mathbf{W}_1 + \mathbf{B}_1^{-1} \sigma_{\alpha_1}^{-2}$ with
 170 $\mathbf{M}_1 = \left(\mathbf{R}_1^{-1} - \mathbf{R}_1^{-1} \mathbf{X}_1 \left(\mathbf{X}'_1 \mathbf{R}_1^{-1} \mathbf{X}_1\right)^{-1} \mathbf{X}'_1 \mathbf{R}_1^{-1}\right)$. Note that only the individual-level dataset 2 and
 171 summary statistics from the dataset 1 (i.e., the estimated allele substitution effects and their
 172 PEC) are required. Individual-level dataset 1 is therefore not required.

173 It is worth noting that the integration of estimates of allele substitution effects from the
 174 dataset 1 into the analysis of dataset 2 can also be obtained from a Bayesian context. Bayes
 175 estimators for linear mixed models were discussed by several authors (Lindley and Smith, 1972;
 176 Dempfle, 1977; Gianola and Fernando, 1986). In a Bayesian context, we can assume the
 177 following prior multivariate normal distributions for the marker model (1) applied to dataset 2:

178 $[\boldsymbol{\beta}_2^* | \mathbf{b}_2, \mathbf{U}_2] \sim MVN(\mathbf{b}_2, \mathbf{U}_2)$, where \mathbf{b}_2 is a mean vector and \mathbf{U}_2 is a (co)variance
 179 matrix,

180 $[\boldsymbol{\alpha}_2^* | \mathbf{B}_2 \sigma_{\alpha_2}^2] \sim MVN(\mathbf{0}, \mathbf{B}_2 \sigma_{\alpha_2}^2)$, and

181 $[\mathbf{e}_2^* | \mathbf{R}_2 \sigma_e^2] \sim MVN(\mathbf{0}, \mathbf{R}_2 \sigma_e^2)$.

182 Assuming a noninformative prior for $\boldsymbol{\beta}_2^*$, the system of equations (2) for dataset 2 can be
 183 obtained by differentiating the joint posterior distribution of $\boldsymbol{\beta}_2^*$ and $\boldsymbol{\alpha}_2^*$ with respect to $\boldsymbol{\beta}_2^*$ and
 184 $\boldsymbol{\alpha}_2^*$, and setting the derivatives equal to 0 (Gianola and Fernando, 1986). Integration of estimates
 185 of allele substitution effects from dataset 1 into the analysis of dataset 2 can be therefore obtained
 186 by defining a multivariate normal prior distribution for allele substitution effects in the analysis
 187 of dataset 2 using the posterior distribution for allele substitution effects from a separate
 188 analysis of dataset 1:

189 $[\boldsymbol{\alpha} | \widehat{\boldsymbol{\alpha}}_1^*, PEC(\widehat{\boldsymbol{\alpha}}_1^*), \mathbf{B}_1 \sigma_{\alpha_1}^2, \mathbf{B}_J \sigma_{\alpha_J}^2] \sim MVN\left(\mathbf{Q} \left(PEC(\widehat{\boldsymbol{\alpha}}_1^*)\right)^{-1} \widehat{\boldsymbol{\alpha}}_1^*, \mathbf{Q}\right)$, (6)

190 $\mathbf{Q} = \left(\left(PEC(\widehat{\boldsymbol{\alpha}}_1^*)\right)^{-1} - \mathbf{B}_1^{-1} \sigma_{\alpha_1}^{-2} + \mathbf{B}_J^{-1} \sigma_{\alpha_J}^{-2} \right)^{-1}$.

191 The matrix \mathbf{Q} can be considered as the PEC matrix of a hypothetical separate analysis of dataset
 192 1 using the multivariate normal prior distribution for allele substitution effects of the joint
 193 analysis, that is $\boldsymbol{\alpha}_1^* \sim MVN(\mathbf{0}, \mathbf{B}_J \sigma_{\alpha_J}^2)$ and $\mathbf{Q} = \left(\mathbf{W}_1' \mathbf{Z}_1' \mathbf{M}_1 \sigma_e^{-2} \mathbf{Z}_1 \mathbf{W}_1 + \mathbf{B}_J^{-1} \sigma_{\alpha_J}^{-2} \right)^{-1}$, and the
 194 vector $\mathbf{Q} \left(PEC(\widehat{\boldsymbol{\alpha}}_1^*)\right)^{-1} \widehat{\boldsymbol{\alpha}}_1^*$ can be considered as the estimated allele substitution effects of this
 195 hypothetical separate analysis. In animal breeding, a similar approach was used to integrate
 196 estimated genetic values and associated accuracies from one genetic evaluation into another
 197 genetic evaluation (Quaas and Zhang, 2006; Legarra et al., 2007; Vandenplas and Gengler,
 198 2012).

199 Finally, it is worth noting that the term $\left(PEC(\widehat{\boldsymbol{\alpha}}_1^*)\right)^{-1} \widehat{\boldsymbol{\alpha}}_1^*$ can be interpreted as a vector
 200 of hypothetical or pseudo-phenotype records associated with allele substitution effects and as

201 such summarize available information in dataset 1. In this sense, the system (5) is similar to
 202 approaches that compute pseudo-records associated with individuals, from available estimated
 203 genetic values where individual-level phenotypic information is not readily available, or is not
 204 measured on the individuals themselves but on close relatives. In animal breeding, these
 205 approaches are commonly known as deregression of estimated genetic values (Jairath et al.,
 206 1998).

207 *Approximate integration*

208 Exact integration requires the inverse of prediction error covariance matrix from the
 209 separate analysis, which could be approximated when unavailable. Genomic analyses of
 210 complex traits that combine different datasets commonly have access to estimated allele
 211 substitution effects and associated prediction error variances (in different forms), but not the
 212 whole prediction error covariance matrix $PEC(\widehat{\boldsymbol{\alpha}}_1^*)$ required in (5). We propose several ways
 213 to accommodate this situation. We assume that we know, at least, the prediction error variances
 214 (PEV) of estimated allele substitution effects ($PEV(\widehat{\boldsymbol{\alpha}}_1^*)$), the number of individuals ($n_{ind,1}$)
 215 and variance components used in the separate analysis of dataset 1 ($\sigma_{\alpha_1}^2$ and σ_e^2).

216 When only the prediction error variances of the estimated allele substitution effects
 217 ($PEV(\widehat{\boldsymbol{\alpha}}_1^*)$) are known, while PEC are not, then we can approximate $(PEC(\widehat{\boldsymbol{\alpha}}_1^*))^{-1}$ with
 218 $(PEV(\widehat{\boldsymbol{\alpha}}_1^*))^{-1}$. This approximation would be accurate if the matrix product $\mathbf{W}'_1\mathbf{W}_1$ has (close
 219 to) zero off-diagonal elements, which is dependent on the characteristics of genotypes in dataset
 220 1 (e.g., allele frequencies, linkage disequilibrium (LD), and population/family structure). If this
 221 is not the case, the approximation will bias the analysis by ignoring off-diagonal elements.

222 When allele frequencies and LD correlations in dataset 1 are known, we can obtain a
 223 good approximation of $PEC(\widehat{\boldsymbol{\alpha}}_1^*)$ under some conditions (one phenotype record per individual,

224 homogenous residual variance, overall mean is the only fixed effect, and Hardy-Weinberg
 225 equilibrium). Derivations in Appendix A2 show that under these conditions we can approximate
 226 $PEC(\widehat{\alpha}_1^*)$ with $(\mathbf{W}'_1\mathbf{W}_1\sigma_e^{-2} + \mathbf{B}_1^{-1}\sigma_{\alpha_1}^{-2})^{-1}$ with the unknown matrix $\mathbf{W}'_1\mathbf{W}_1$ approximated
 227 from commonly available population parameters (i.e., allele frequencies and LD correlation) as
 228 $4n_{ind,1}\mathbf{pp}' + \mathbf{V}^{\frac{1}{2}}\mathbf{CV}^{\frac{1}{2}}$, where \mathbf{p} is a $n_{mar} \times 1$ vector of allele frequencies, \mathbf{V} is a $n_{mar} \times n_{mar}$
 229 diagonal matrix of expected genotype sum of squares with the i -th diagonal element equal to
 230 $n_{ind,1}2p_{i,1}(1 - p_{i,1})$, and \mathbf{C} is a $n_{mar} \times n_{mar}$ matrix of pairwise genotype correlations between
 231 markers. In practice, the matrix \mathbf{C} for dataset 1 could be unknown, but we can approximate it
 232 by using a reference panel that includes, for example, available genotypes of non-phenotyped
 233 individuals originating from this population (Yang et al., 2012; Vilhjálmsson et al., 2015; Maier
 234 et al., 2018).

235 Finally, we relax the assumption of having a single phenotype record per individual in
 236 the preceding approximations. This is relevant when individuals have repeated phenotype
 237 records, e.g., repeated longitudinal production or reproduction records in livestock or replicated
 238 field trials in crops. A related issue is the violation of assumption of homogenous residual
 239 variance when phenotype records are first pre-processed and then used in genomic analyses,
 240 e.g., deregressed progeny proofs in livestock (e.g., Garrick et al., 2009) or adjusted field trial
 241 means in crops (e.g., Schulz-Streeck et al., 2013; Oakey et al., 2016; Damesa et al., 2017). For
 242 these situations, we show in Appendix A3 that we can approximate $PEC(\widehat{\alpha}_1^*)$ with
 243 $(\Lambda_1(4\mathbf{pp}' + \Psi^{\frac{1}{2}}\mathbf{C}\Psi^{\frac{1}{2}})\Lambda_1\sigma_e^{-2} + \mathbf{B}_1^{-1}\sigma_{\alpha_1}^{-2})^{-1}$ where Ψ is a $n_{mar} \times n_{mar}$ diagonal matrix with
 244 the j -th diagonal element equal to $2p_{j,1}(1 - p_{j,1})$, and Λ_1 is a $n_{mar} \times n_{mar}$ diagonal matrix
 245 with the j -th diagonal element representing the square root of effective number of records for
 246 the j -th marker. The matrix Λ_1 can be obtained by solving the nonlinear system of equations

$$247 \quad \text{diag} \left(\left(\mathbf{\Lambda}_1 \left(4\mathbf{p}\mathbf{p}' + \mathbf{\Psi}^{\frac{1}{2}}\mathbf{C}\mathbf{\Psi}^{\frac{1}{2}} \right) \mathbf{\Lambda}_1 \sigma_e^{-2} + \mathbf{B}_1^{-1} \sigma_{\alpha_1}^{-2} \right)^{-1} \right) = PEV(\widehat{\boldsymbol{\alpha}}_1^*)$$

248 through a fixed-point iteration algorithm (Burden and Faires, 2010) detailed in Appendix A3.

249 It is worth noting that the proposed algorithm requires the inversion of a $n_{mar} \times n_{mar}$ dense
 250 matrix at each iteration. This computational cost can be reduced by performing the algorithm
 251 for each chromosome separately.

252 *Integration with multiple populations*

253 When more than two populations or datasets are available the developed methods can
 254 be easily extended. With n datasets, the prior distribution for allele substitution effects in the
 255 separate analysis of the n -th dataset is defined using the posterior distributions for allele
 256 substitution effects from the separate analyses of $n - 1$ datasets:

$$257 \quad [\boldsymbol{\alpha} | \widehat{\boldsymbol{\alpha}}_1^*, \widehat{\boldsymbol{\alpha}}_2^*, \dots, \widehat{\boldsymbol{\alpha}}_{n-1}^*] \sim MVN \left(\mathbf{Q} \sum_{i=1}^{n-1} \left((PEC(\widehat{\boldsymbol{\alpha}}_i^*))^{-1} \widehat{\boldsymbol{\alpha}}_i^* \right), \mathbf{Q} \right),$$

$$258 \quad \mathbf{Q} = \left(\mathbf{B}_J^{-1} \sigma_{\alpha_j}^{-2} + \sum_{i=1}^{n-1} \left((PEC(\widehat{\boldsymbol{\alpha}}_i^*))^{-1} - \mathbf{B}_i^{-1} \sigma_{\alpha_i}^{-2} \right) \right)^{-1}.$$

259 **Simulations**

260 We tested developed methods with simulated data that either had low or high genetic
 261 diversity. The data was simulated in 5 replicates with the AlphaSim program, which uses the
 262 coalescent method for simulation of base population chromosomes and the gene drop method
 263 for simulation of chromosome inheritance within a pedigree (Hickey and Gorjanc, 2012; Faux
 264 et al., 2016).

265 A diploid genome was simulated with 30 chromosomes, each 10^8 base pairs long.
 266 Coalescent mutation and recombination rate per base pair were set to 10^{-8} , while effective
 267 population size was modelled over time to mimic population history of a livestock population

Multi-population genomic prediction

268 in line with the values reported by MacLeod et al. (2013). Specifically, for the low diversity
269 scenario effective population size of the base population was set to 100 and increased to 120,
270 250, 350, 1,000, 1,500, 2,000, 2,500, 3,500, 7,000, 10,000, 17,000, and 62,000 at respectively
271 6, 12, 18, 24, 154, 454, 654, 1,754, 2,354, 3,354, 33,154, and 933,154 generations ago. For the
272 high diversity scenario, effective population size of the base population was set to 10,000 and
273 increased above this value in the same way as in the low diversity scenario; to 17,000 and
274 62,000 at 33,154, and 933,154 generations ago. For each chromosome 10,000 whole
275 chromosome haplotypes were sampled, which on average hosted about 700,000 markers (21
276 million per genome) for the low diversity scenario and 1,400,000 markers (42 million per
277 genome) for the high diversity scenario. Out of these loci 100 per chromosome (3,000 per
278 genome) were sampled as causal loci affecting a complex trait. The allele substitution effect of
279 causal loci was sampled from a normal distribution with mean zero and variance $1/3,000$. The
280 effects were used to simulate a complex trait with additive genetic architecture. In addition,
281 2,000 loci per chromosome (60,000 per genome) were selected as markers with the restriction
282 of having minor allele frequency above 0.05.

283 From the base population, founder genomes for four populations (A, B, C, and D) were
284 obtained by random sampling of chromosomes with recombination. The populations were
285 ancestrally related through the common base population, but otherwise maintained
286 independently, i.e., there was no migration between the four populations. Each population was
287 initiated with 10,000 founders (half males and half females) and maintained for 7 generations
288 with constant size. In the low diversity scenario, with the effective population size of 100, 25
289 males and 5,000 females were selected as parents of each generation, while in the high diversity
290 scenario, with the effective population size of 10,000, all 5,000 males and 5,000 females were
291 used. The 25 males were selected on true genetic value, assuming accurate progeny test was
292 available.

Multi-population genomic prediction

293 For every individual in the population we simulated two types of phenotypes. First, an
294 own single phenotype was simulated as the sum of the true genetic value and a residual sampled
295 from a normal distribution with mean zero and residual variance scaled relative to the variance
296 of true genetic value in the base population such that heritability was 0.3. These simulated single
297 phenotype records mimic records measured on the individual. Second, a weighted phenotype
298 was simulated as the sum of the true genetic value and the mean of n_{weight} residuals. Each
299 residual was sampled from a normal distribution with mean zero and residual variance scaled
300 relative to the variance of true genetic value in the base population such that heritability was
301 0.3. The weight n_{weight} was equal to $n_{weight} = 1 + val$ where the real value val was sampled
302 from a geometric distribution with a probability p of 0.15 and a probability mass function of
303 $Pr(x) = p(1 - p)^x$ with $x \in \{0, 1, 2, \dots\}$. The average n_{weight} was 6.6. These weighted
304 phenotypes mimic either repeated records of an individual or records on multiple progeny of an
305 individual. To satisfy the assumption of identical residual variance across all analyses,
306 phenotype records were divided by the residual standard deviation specific for each population,
307 such that $\sigma_e^2 = 1$. For every individual in each population we stored the true genetic value, own
308 single and weighted phenotype records, associated weight, and 60,000 marker genotypes.

309 **Analysis**

310 The data was analysed in several ways to evaluate the developed methods. In each case
311 the aim was to obtain accurate genetic values utilizing all the available information.
312 Specifically, we integrated results from separate analysis of populations B, C, and D, into the
313 analysis of population A. We assumed throughout that variance components were known and
314 equal to the rescaled variances. We analysed three scenarios in total. The first and second
315 scenario used population specific training data of randomly sampled 30,000 individuals with
316 single phenotype record from generations 1 to 6 under low and high diversity settings. The third

Multi-population genomic prediction

317 scenario used population specific training data of randomly sampled 10,000 individuals with
318 weighted phenotype record from generations 1 to 6 under low diversity setting. In all scenarios
319 all of the 10,000 individuals from generation 7 of each population were considered as validation
320 individuals. The following analyses were performed:

- 321 1) A joint analysis of four populations. This was the reference that the other analyses
322 were compared against;
- 323 2) A separate analysis for each of the four populations;
- 324 3) An exact integration of separate analyses of populations B, C, and D, into the
325 analysis of population A;
- 326 4) The same as 3), but approximating the PEC matrix with a partial PEC matrix for
327 each chromosome, i.e., PEC between markers on different chromosomes were set
328 to zero;
- 329 5) The same as 3), but approximating the PEC matrix with a diagonal PEV matrix, i.e.,
330 PEC between all markers were set to zero;
- 331 6) The same as 3), but approximating the PEC matrix with PEV, allele frequencies,
332 and LD correlations between markers obtained from the training sets. For the
333 scenario with weighted phenotype records, the algorithm for estimating the effective
334 number of records per marker was performed for each marker separately and for
335 each chromosome separately.
- 336 7) The same as 6), but with LD correlations between markers computed from
337 validation individuals instead of the training data.

338 For each analysis we calculated genomic prediction accuracy as the Pearson correlation
339 between the true and estimated genetic value in validation individuals. Further, we evaluated
340 the different integrations by comparing estimated genetic values of validation individuals
341 against the estimated genetic values obtained from the joint analysis, which was considered as

Multi-population genomic prediction

342 the reference because it used information from all populations. If integration was fully accurate,
343 there should be no difference between the joint analysis and the analysis with integration. We
344 assessed this by (a) accuracy of integration as a Pearson correlation between estimated genetic
345 values from the joint analysis and the analysis with integration (desired value equals 1), (b)
346 calibration of integration as a regression of estimated genetic values from the joint analysis on
347 estimated genetic values from analysis with integration, and (c) magnitude of error in
348 integration as a mean square error (MSE) between estimated genetic values from the joint
349 analysis and from the analysis with integration (desired value equals 0). By calibration, we
350 mean the slope of relationship of the estimates from the integration analysis onto the estimated
351 genetic values from the joint analysis. The desired slope value is 1, which indicates a well
352 calibrated model. Values above or below 1 indicate an uncalibrated model.

353 **Data availability**

354 Supplemental figures are available in File S1. A description of the simulated genotype
355 and phenotype datasets for each scenario is provided in File S2. Simulated genotype and
356 phenotype datasets for the 5 replicates of each scenario are provided in Files S3, S4, and S5.
357 Data simulation scripts and Fortran codes developed to perform the different analyses, as well
358 as a short description of each of them, are provided in File S6. All files were uploaded to
359 Figshare: <https://figshare.com/s/473dc83a7b154cfd610c>.

360

361

RESULTS

362 **Genomic prediction accuracy of separate and joint analyses**

363 Joint analysis increased genomic prediction accuracy in comparison to separate
364 analyses. This is shown in Table 1. Analysing separately the four datasets gave accuracies of
365 about 0.71 (low diversity) and 0.53 (high diversity) with single phenotype records, and of about
366 0.73 (low diversity) with weighted phenotype records. Analysing jointly the four datasets
367 increased accuracy by at least 0.09 absolute points with single phenotype records and by at least
368 0.12 absolute points with weighted phenotype records.

369 **Integration based on PEC, partial PEC, or PEV matrices**

370 For all scenarios the developed method enabled exact integration when complete PEC
371 matrices were used. Integration of estimated allele substitution effects by means of the complete
372 PEC matrix led to the same estimated genetic values as with the joint analysis, as shown by
373 correlation and regression coefficients of 1, and MSE close to 0 (Figures 1-4; Figures S1-S8).
374 For comparison, correlations between estimated genetic values from separate analyses and joint
375 estimated genetic values were about 0.87 (low diversity) and 0.77 (high diversity) with single
376 phenotype records, and 0.85 (low diversity) with weighted phenotype records.

377 Approximate integration by means of partial PEC matrices for each chromosome, that
378 is ignoring PEC between markers on different chromosomes, gave almost as accurate and
379 calibrated estimated genetic values as the exact integration. This is illustrated in Figures 1-4
380 and Figures S1-S8 with correlations higher than 0.96, regression coefficients close to 1, and
381 MSE close to 0. Increasing the diversity slightly deteriorated accuracy and calibration of
382 genomic predictions (Figures 1-2; Figures S1-S4).

383 Approximate integrations by means of PEV matrices, that is ignoring PEC between all
384 markers, gave quite accurate, but not calibrated estimated genetic values. This is shown in
385 Figures 1-4 and in Figures S1-S8. Correlations between joint estimated genetic values and
386 estimated genetic values with integration by means of PEV were between 0.95 and 0.98 with
387 single phenotype records and between 0.93 and 0.95 with weighted phenotype records. Despite
388 these correlations close to 1, estimated genetic values were not well calibrated, as depicted by
389 regression coefficients below 0.77 for the low diversity scenarios with single and weighted
390 phenotype records, and below 0.86 for the high diversity scenario with single phenotype records
391 (Figures 2, 4, S2, S6).

392 **Integration based on PEV, allele frequencies, and LD information**

393 When LD information was derived from training data of other populations, approximate
394 integrations by means of PEV, allele frequencies, and LD information, resulted in highly
395 accurate and well calibrated estimated genetic values with single phenotype records. This is
396 shown in Figures 1-2 (Figures S1-S4). Correlation and regression coefficients were equal to 1
397 for the low diversity scenario. Slightly lower values, but still close to 1, were observed for the
398 high diversity scenario. For both low and high diversity scenarios, MSE were close to 0. In
399 contrast, when LD information was derived from validation data of other populations,
400 approximate integrations gave less accurate and calibrated estimated genetic values. This is
401 shown in Figures 1-2 (Figures S1-S4). For these scenarios, correlations were equal to at least
402 0.94, and regression coefficients varied between 0.87 and 1.05.

403 For the scenario with weighted phenotype records, approximate integrations by means
404 of LD information from training data of other populations resulted in highly accurate and well
405 calibrated estimated genetic values when sets of markers per chromosome were used to estimate
406 the effective number of records for each marker. Correlations between joint estimated genetic

Multi-population genomic prediction

407 values and estimated genetic values with integration were about 0.99 (Figure 3, Figure S5),
408 regression coefficients were about 0.95 (Figure 4, Figure S6), and MSE were close to 0 (Figures
409 S7-S8). Using LD information from the validation data of other populations, instead from the
410 training data of other populations, gave slightly less accurate (correlations higher than 0.95),
411 and moderately less calibrated estimated genetic values (regression coefficients between 0.87
412 and 1.04; Figures 3-4; Figures S5-S8). For both cases, estimating the effective numbers of
413 records per marker, instead of for all markers per chromosome simultaneously, reduced
414 accuracy and calibration of estimated genetic values (Figures 3-4; Figures S5-S6).

415 **Comparison of estimated allele substitution effects**

416 Correlation and regression coefficients between estimated allele substitution effects
417 from the joint analysis and analysis with integration largely followed patterns of the
418 corresponding values for estimated genetic values (Tables 2-3). Correlation and regression
419 coefficients were close to 1 when the integration of estimated allele substitution effects was by
420 means of the complete PEC matrices. Ignoring PEC between markers on different
421 chromosomes, or ignoring PEC between all markers, reduced correlations to between 0.92 and
422 0.99 (Tables 2-3). Using LD information with PEV led to correlations between joint estimates
423 of allele substitution effects and estimates with integration ranging from 0.71 to 0.83 for the
424 scenario with weighted phenotype records (Tables 2-3).

425

426

DISCUSSION

427 The results show that the developed method enables accurate and well calibrated
428 estimated genetic values for complex traits using both individual-level data and summary
429 statistics. As expected from theory, the analysis of individual-level data and estimated allele
430 substitution effects from other analyses by means of PEC matrices, yielded the same estimates
431 as the joint analysis of all individual-level data. To our knowledge, this is the first time that
432 individual-level data and summary statistics were analysed simultaneously for genomic
433 predictions. As illustrated by simulations, the combined analysis of multiple datasets may
434 increase genomic prediction accuracy over separate analyses of a single dataset. Unfortunately,
435 combining individual-level data from several sources is generally not feasible for several
436 reasons, e.g., political roadblocks, data protection concerns, or data inconsistencies (Powell
437 and Sieber, 1992; Vilhjálmsson et al., 2015; Maier et al., 2018). However, summary statistics,
438 such as estimates of allele substitution effects and associated measures of accuracy (e.g., PEV),
439 are usually available for exchange in human genetics, or are discussed to be shared, e.g., at an
440 international level for dairy cattle breeding (Liu and Goddard, 2018). The developed method
441 enables increase in genomic prediction accuracy of complex traits by means of jointly analysing
442 the available individual-level data and summary statistics.

443 Accurate integration of estimated allele substitution effects is possible also when the
444 complete PEC matrix is not available. This is important because computing the exact PEC
445 matrix and exchanging it between analyses might be challenging in some cases. For the vast
446 majority of used marker arrays in animal and plant breeding the calculations and data transfers
447 should be doable. For example, most arrays have between 10,000 and 100,000 markers, for
448 which we need between ~1 and ~80 GB of memory to store the PEC matrix and between a
449 minute and a day to invert it on current computers. For a larger number of markers, commonly
450 used in human genetics, the memory requirements and computing time become prohibitive. The

Multi-population genomic prediction

451 results show that in such cases we can still obtain accurate genomic predictions when the
452 integration is done by means of partial PEC matrices for each chromosome. This is expected
453 since high LD between markers mostly occurs within chromosomes. High LD between markers
454 on different chromosomes may especially occur in structured populations and populations
455 under selection (Farnir et al., 2000; Flint-Garcia et al., 2003; Rostoks et al., 2006). Both of these
456 conditions are present in breeding populations. However, the results suggest that LD between
457 chromosomes can be ignored for the purpose of integration for populations with both low and
458 high diversity. The results also show that we can successfully integrate estimated allele
459 substitution effects when only PEV and allele frequencies from each population are available
460 together with LD information of a reference genotype panel representative of each population.
461 Assuming that such reference genotype panels are available, only estimated allele substitution
462 effects, associated PEV, and allele frequencies need to be exchanged between populations for
463 such analyses. Similar conclusions were drawn from studies combining only summary statistics
464 obtained from genome-wide association studies to perform multi-trait genomic predictions
465 (Maier et al., 2018).

466 Accurate integration of estimated allele substitution effects is possible irrespective of
467 the diversity of the populations and characteristics of genotypes (e.g., allele frequencies, LD).
468 This is obvious, and confirmed by our results, when integration is performed by means of
469 complete PEC matrices. When complete PEC matrices are unavailable, accurate integration is
470 possible if the inverses of the PEC matrices can be approximated accurately from available
471 population parameters (i.e. LD and allele frequency information), whatever the level of
472 diversity and characteristics of the populations, as shown by our results or a study combining
473 summary statistics in human genetics (Maier et al., 2018). In our study, the population
474 parameters obtained from the reference panels adequately reflected the characteristics of the
475 training sets. We expect that this would be the case for populations with substantial migration,

Multi-population genomic prediction

476 such as, for example, Holstein dairy cattle populations. Future studies should be conducted to
477 assess the impact of suboptimal reference panels. Therefore, the developed method is expected
478 to perform well on any type of data, from animal and plant breeding to human genetics,
479 provided accurate information is available.

480 The developed method has some simplifying assumptions that can be readily relaxed.
481 For example, we assumed that the same genotype coding was used in all populations. This
482 assumption can be relaxed when centered genotype coding (i.e., of the form of $(\mathbf{W}_i - \mathbf{1}\mathbf{v}_i')$) is
483 used because variance component estimates, estimates of allele substitution effects and PEC
484 are the same irrespective of the centering of the genotype coding, provided that the model has
485 a fixed general mean, which is considered in the integration (Strandén and Christensen, 2011).
486 Also, centered and scaled (standardised) genotype coding is often used in human genetics,
487 instead of only centered genotype coding (Yang et al., 2010; Speed et al., 2012; Maier et al.,
488 2018). In practice, estimates of genetic values are only slightly influenced by scaling of centered
489 genotype coding (Strandén and Christensen, 2011; Bouwman et al., 2017). Therefore, assuming
490 that the same estimated genetic values are obtained with different scaling, allele substitution
491 effects estimated using one type of genotype scaling could be obtained from a post-analysis by
492 converting estimated genetic values computed for a reference genotype panel into allele
493 substitution effects for another genotype scaling. Converting estimated genetic values into
494 allele substitution effects is often referred to as back-solving of allele substitution effects
495 (Strandén and Garrick, 2009; Strandén and Christensen, 2011; Wang et al., 2012; Bouwman et
496 al., 2017). Prediction error covariances associated with the converted estimated allele
497 substitution effects could be derived from the (prediction error) covariances of the estimated
498 genetic values (see derivations in Appendix A4).

499 Allele substitution effects estimated from analyses using different sets of markers or
500 different residual variances, can be used in the integration as well. The assumption that all

501 individuals were genotyped at the same loci could be considered as fulfilled if small differences
 502 in the sets of markers are corrected by assuming zero allele substitution effect and zero accuracy
 503 for markers not used in an analysis. When large differences between sets of markers are
 504 observed, this assumption can be accomodated following two approaches. A first, post-analysis,
 505 approach consists of assuming that estimated genetic values are the same for two different sets
 506 of markers, allowing the conversion of estimated allele substitution effects from one set of
 507 markers to another set of markers (Liu and Goddard, 2018). The conversion can be performed
 508 by back-solving estimated allele substitution effects from estimated genetic values, as proposed
 509 previously for different genotype codings, or by applying a marker model to the estimated
 510 genetic values with the reference set of markers (Liu and Goddard, 2018). A second approach
 511 consists of harmonizing genotype data across populations. This approach must be performed
 512 before the analyses, and requires therefore coordination between populations. Harmonization
 513 of genotype data could be performed by identifying a subset of markers for which all
 514 populations are genotyped, or by genotype imputation (e.g., Marchini and Howie, 2010).
 515 Finally, the assumption that residual variances were the same in all populations, can be relaxed
 516 by noting that separate estimates of allele substitution effects $\widehat{\alpha}_i^*$, obtained by the system of
 517 equations (2), can be also obtained by the following different formulations:

$$\begin{aligned}
 \widehat{\alpha}_i^* &= (\mathbf{W}_i' \mathbf{Z}_i' \mathbf{M}_i \sigma_{e_i}^2 \mathbf{Z}_i \mathbf{W}_i + \mathbf{B}_i^{-1} \sigma_{\alpha_i}^{-2})^{-1} \mathbf{W}_i' \mathbf{Z}_i' \mathbf{M}_i \sigma_{e_i}^2 \mathbf{y}_i \\
 &= (\mathbf{W}_i' \mathbf{Z}_i' \mathbf{M}_i \mathbf{Z}_i \mathbf{W}_i + \mathbf{B}_i^{-1} \lambda)^{-1} \mathbf{W}_i' \mathbf{Z}_i' \mathbf{M}_i \mathbf{y}_i \\
 &= (\mathbf{W}_i' \mathbf{Z}_i' \mathbf{M}_i \sigma_{e_f}^{-2} \mathbf{Z}_i \mathbf{W}_i + \mathbf{B}_i^{-1} \lambda \sigma_{e_f}^{-2})^{-1} \mathbf{W}_i' \mathbf{Z}_i' \mathbf{M}_i \sigma_{e_f}^{-2} \mathbf{y}_i
 \end{aligned}$$

519 where $\sigma_{e_i}^2$ ($\sigma_{e_f}^2$) is the residual variance used for the i -th (focal) analysis, and $\lambda = \sigma_{e_i}^2 \sigma_{\alpha_i}^{-2}$.

520 For integration of $\widehat{\alpha}_i^*$, $(PEC(\widehat{\alpha}_i^*))^{-1}$ must be approximated using the residual variance of the
 521 focal population ($\sigma_{e_f}^2$) and the effective numbers of records per marker estimated using variance
 522 components of the i -th analysis. Another way to relax this assumption is to extend our univariate

523 model to a bivariate model, similarly to methods developed to combine different genetic
524 evaluations in animal breeding (Schaeffer, 1994; Vandenplas et al., 2015). In a bivariate model,
525 one trait would represent individual-level data, while the other trait would represent summary
526 statistics. The genetic correlation between the two traits could be estimated based on a subset
527 of individual-level data available for both datasets or based on summary statistics (Bulik-
528 Sullivan et al., 2015). Such an approach would also allow the integration of summary statistics
529 expressed on a different scale (e.g., different measure units, trait definitions) than the scale of
530 the focal population (Vandenplas et al., 2015).

531 The developed method can be readily generalized to multi-trait models and is therefore
532 a generalization of previous works that were based on several (implicit) assumptions (Liu and
533 Goddard, 2018; Maier et al., 2018). For example, previous works assumed that no individual-
534 level data were available. It was also (implicitly) assumed that only single phenotype records
535 with homogeneous residual variance (Maier et al., 2018), or that the least-squares part of the
536 separate analyses (Liu and Goddard, 2018), were available for integrating estimated allele
537 substitution effects. Both assumptions lead to simple and accurate approximations of PEC
538 matrices as shown in our study. However, we relax all these assumptions, such that our method
539 can jointly analyse individual-level data and summary statistics, with possibly multiple
540 phenotype records per individual.

541 With all the proposed generalizations, the developed method could be used in different
542 contexts. For example, in human genetics, allele substitution effects with associated standard
543 errors are publicly available (Yang et al., 2012; Vilhjálmsson et al., 2015; Maier et al., 2018).
544 In animal breeding, individuals' genetic values with associated reliabilities are publicly
545 available and in the case of dairy cattle extensively combined across multiple populations
546 (Schaeffer, 1994; VanRaden and Sullivan, 2010; Jorjani et al., 2012; Vandenplas et al., 2017).
547 The developed method can be used in both contexts, but in the latter case individuals' genetic

Multi-population genomic prediction

548 values must be first back-solved to allele substitution effects (Strandén and Garrick, 2009;
549 Strandén and Christensen, 2011; Wang et al., 2012; Bouwman et al., 2017). It is worth noting
550 that our method assumes that summary statistics from one population are free of information
551 from other populations. This suggest that it can be used when there is no, or limited, sharing of
552 information between populations, as is for example the case in beef cattle, but not in dairy cattle
553 populations such as Holstein, where pseudo-phenotypes summarising information from
554 multiple populations are used extensively (VanRaden and Sullivan, 2010; Jorjani et al., 2012).
555 This assumption can be relaxed by performing separate analyzes free of information from other
556 populations, or by correcting for double-counting of information, which has bee developed for
557 the integration of estimated genetic values from different populations (Vandenplas et al., 2014,
558 2017; VanRaden et al., 2014). This correction for double-counting of information is not yet
559 developed for the integration of summary statistics, and should be investigated in future studies.

560

561

CONCLUSIONS

562 We developed a method for genomic prediction that accurately integrates summary
563 statistics obtained from analyses of separate populations into an analysis of individual-level
564 data. The method accommodates use of multiple phenotype (pseudo-)records per individual,
565 and further extensions have been presented to accommodate for differences in residual
566 variances or genotype codings used in the populations. When complete summary statistics
567 information is available the method gives identical genomic predictions as the joint analysis of
568 individual-level data from all populations. When summary statistics information is not
569 complete we can use a series of approximations that give very accurate and well calibrated
570 genomic predictions.

571

572

AUTHORS' CONTRIBUTIONS

573

JV derived the equations, wrote the programs to do the analyses, performed the

574

analyses, and drafted the outline of the manuscript. GG performed the simulations. All

575

authors discussed the design of the simulations. JV and GG wrote the first version of the

576

manuscript. All authors provided valuable insights throughout the analysis and writing

577

process.

578

579

ACKNOWLEDGMENTS

580 This study was financially supported by the Dutch Ministry of Economic Affairs (TKI
581 Agri & Food project 16022), the Breed4Food partners Cobb Europe, CRV, Hendrix Genetics
582 and Topigs Norsvin, and UK Biotechnology and Biological Sciences Research Council
583 (BBSRC) ISPG to The Roslin Institute BBS/E/D/30002275. The use of the HPC cluster has
584 been made possible by CAT-AgroFood (Shared Research Facilities Wageningen UR).

585

586

LITERATURE CITED

- 587 Bouwman, A.C., B.J. Hayes, and M.P.L. Calus. 2017. Estimated allele substitution effects
588 underlying genomic evaluation models depend on the scaling of allele counts. *Genet.*
589 *Sel. Evol.* 49. doi:10.1186/s12711-017-0355-9.
- 590 Bulik-Sullivan, B., H.K. Finucane, V. Anttila, A. Gusev, F.R. Day, et al. 2015. An atlas of
591 genetic correlations across human diseases and traits. *Nat. Genet.* 47:1236–1241.
592 doi:10.1038/ng.3406.
- 593 Burden, R.L., and J.D. Faires. 2010. *Numerical Analysis*. 9 edition. Brooks Cole, Boston,
594 MA.
- 595 Campos, G. de los, D. Gianola, and D.B. Allison. 2010. Predicting genetic predisposition in
596 humans: the promise of whole-genome markers. *Nat. Rev. Genet.* 11:880–886.
597 doi:10.1038/nrg2898.
- 598 de los Campos, G., J.M. Hickey, R. Pong-Wong, H.D. Daetwyler, and M.P.L. Calus. 2012.
599 Whole-genome regression and prediction methods applied to plant and animal
600 breeding. *Genetics* 193:327–345. doi:10.1534/genetics.112.143313.
- 601 Daetwyler, H.D., B. Villanueva, and J.A. Woolliams. 2008. Accuracy of predicting the
602 genetic risk of disease using a genome-wide approach. *PLoS ONE* 3.
- 603 Damesa, T.M., J. Möhring, M. Worku, and H.-P. Piepho. 2017. One step at a time: Stage-wise
604 analysis of a series of experiments. *Agron. J.* 109:845–857.
605 doi:10.2134/agronj2016.07.0395.
- 606 Dempfle, L. 1977. Relation entre BLUP (Best Linear Unbiased Prediction) et estimateurs
607 bayésiens. *Genet. Sel. Evol.* 9:27–32.

Multi-population genomic prediction

- 608 Farnir, F., W. Coppieters, J.-J. Arranz, P. Berzi, N. Cambisano, et al. 2000. Extensive
609 genome-wide linkage disequilibrium in cattle. *Genome Res.* 10:220–227.
610 doi:10.1101/gr.10.2.220.
- 611 Faux, A.-M., G. Gorjanc, R.C. Gaynor, M. Battagin, S.M. Edwards, et al. 2016. AlphaSim:
612 Software for breeding program simulation. *Plant Genome* 9.
- 613 Fisher, R.A. 1918. The correlation between relatives on the supposition of Mendelian
614 inheritance. *Philos. Trans. R. Soc. Edinb.* 52:399–433.
- 615 Flint-Garcia, S.A., J.M. Thornsberry, and E.S. Buckler. 2003. Structure of linkage
616 disequilibrium in plants. *Annu. Rev. Plant Biol.* 54:357–374.
617 doi:10.1146/annurev.arplant.54.031902.134907.
- 618 Garrick, D.J., J.F. Taylor, and R.L. Fernando. 2009. Deregressing estimated breeding values
619 and weighting information for genomic regression analyses. *Genet. Sel. Evol.* 41:55.
620 doi:10.1186/1297-9686-41-55.
- 621 Gianola, D., and R.L. Fernando. 1986. Bayesian methods in animal breeding theory. *J. Anim.*
622 *Sci.* 63:217–244.
- 623 Henderson, C.R. 1984. *Applications of Linear Models in Animal Breeding*. 2nd ed.
624 University of Guelph, Guelph, ON, Canada.
- 625 Hickey, J.M., and G. Gorjanc. 2012. Simulated data for genomic selection and genome-wide
626 association studies using a combination of coalescent and gene drop methods. *G3*
627 2:425–427. doi:10.1534/g3.111.001297.

- 628 Hoerl, A.E., and R.W. Kennard. 1976. Ridge regression iterative estimation of the biasing
629 parameter. *Commun. Stat. - Theory Methods* 5:77–88.
630 doi:10.1080/03610927608827333.
- 631 Hozé, C., S. Fritz, F. Phocas, D. Boichard, V. Ducrocq, et al. 2014. Efficiency of multi-breed
632 genomic selection for dairy cattle breeds with different sizes of reference population.
633 *J. Dairy Sci.* 97:3918–3929. doi:10.3168/jds.2013-7761.
- 634 Jairath, L., J.C.M. Dekkers, L.R. Schaeffer, Z. Liu, E.B. Burnside, et al. 1998. Genetic
635 evaluation for herd life in Canada. *J. Dairy Sci.* 81:550–562.
- 636 Jorjani, H., J. Jakobsen, E. Hjerpe, V. Palucci, and J. Dürr. 2012. Status of genomic
637 evaluation in the Brown Swiss populations. *Interbull Bull.* 46:46–54.
- 638 Legarra, A., J.K. Bertrand, T. Strabel, R.L. Sapp, J.P. Sanchez, et al. 2007. Multi-breed
639 genetic evaluation in a Gelbvieh population. *J. Anim. Breed. Genet.* 124:286–295.
- 640 Lindley, D.V., and A.F.M. Smith. 1972. Bayes estimates for the linear model. *J. R. Stat. Soc.*
641 *Ser. B Methodol.* 34:1–41.
- 642 Liu, Z., and M.E. Goddard. 2018. A SNP MACE model for international genomic evaluation:
643 technical challenges and possible solutions. Page 11.393 in *Proceedings of the 11th*
644 *World Congress on Genetics Applied to Livestock Production*, Auckland, New
645 Zealand.
- 646 MacLeod, I.M., D.M. Larkin, H.A. Lewin, B.J. Hayes, and M.E. Goddard. 2013. Inferring
647 demography from runs of homozygosity in whole-genome sequence, with correction
648 for sequence errors. *Mol. Biol. Evol.* 30:2209–2223.

Multi-population genomic prediction

- 649 Maier, R.M., Z. Zhu, S.H. Lee, M. Trzaskowski, D.M. Ruderfer, et al. 2018. Improving
650 genetic prediction by leveraging genetic correlations among human diseases and traits.
651 Nat. Commun. 9:989.
- 652 Marchini, J., and B. Howie. 2010. Genotype imputation for genome-wide association studies.
653 Nat. Rev. Genet. 11:499–511. doi:10.1038/nrg2796.
- 654 Meuwissen, T.H.E., B.J. Hayes, and M.E. Goddard. 2001. Prediction of total genetic value
655 using genome-wide dense marker maps. Genetics 157:1819–1829.
- 656 Misztal, I., and G.R. Wiggans. 1988. Approximation of prediction error variance in large-
657 scale animal models. J. Dairy Sci. 71(Suppl. 2):27–32.
- 658 Oakey, H., B. Cullis, R. Thompson, J. Comadran, C. Halpin, et al. 2016. Genomic selection in
659 multi-environment crop trials. G3 Bethesda Md 6:1313–1326.
660 doi:10.1534/g3.116.027524.
- 661 Pasaniuc, B., and A.L. Price. 2017. Dissecting the genetics of complex traits using summary
662 association statistics. Nat. Rev. Genet. 18:117–127. doi:10.1038/nrg.2016.142.
- 663 Powell, R.L., and H.D. Norman. 1998. Use of multinational data to improve national
664 evaluations of Holstein bulls. J. Dairy Sci. 81:2257–2263. doi:10.3168/jds.S0022-
665 0302(98)75805-9.
- 666 Powell, R.L., and M. Sieber. 1992. Direct and indirect conversion of bull evaluations for yield
667 traits between countries. J. Dairy Sci. 75:1138–1146.
- 668 Quaas, R.L., and Z. Zhang. 2006. Multiple-breed genetic evaluation in the US beef cattle
669 context: Methodology. Page CD-ROM Comm. 24-12 in Proceedings of the 8th World
670 Congress on Genetics Applied to Livestock Production, Belo Horizonte, Brazil.

Multi-population genomic prediction

- 671 Rogers, A.R., and C. Huff. 2009. Linkage Disequilibrium between loci with unknown phase.
672 *Genetics* 182:839–844. doi:10.1534/genetics.108.093153.
- 673 Rostoks, N., L. Ramsay, K. MacKenzie, L. Cardle, P.R. Bhat, et al. 2006. Recent history of
674 artificial outcrossing facilitates whole-genome association mapping in elite inbred
675 crop varieties. *Proc. Natl. Acad. Sci. U. S. A.* 103:18656–18661.
676 doi:10.1073/pnas.0606133103.
- 677 Schaeffer, L.R. 1994. Multiple-country comparison of dairy sires. *J. Dairy Sci.* 77:2671–2678.
- 678 Schulthess, A.W., Y. Wang, T. Miedaner, P. Wilde, J.C. Reif, et al. 2016. Multiple-trait- and
679 selection indices-genomic predictions for grain yield and protein content in rye for
680 feeding purposes. *TAG Theor. Appl. Genet. Theor. Angew. Genet.* 129:273–287.
681 doi:10.1007/s00122-015-2626-6.
- 682 Schulz-Streeck, T., J.O. Ogutu, and H.-P. Piepho. 2013. Comparisons of single-stage and two-
683 stage approaches to genomic selection. *Theor. Appl. Genet.* 126:69–82.
684 doi:10.1007/s00122-012-1960-1.
- 685 Speed, D., G. Hemani, M.R. Johnson, and D.J. Balding. 2012. Improved heritability
686 estimation from genome-wide SNPs. *Am. J. Hum. Genet.* 91:1011–1021.
687 doi:10.1016/j.ajhg.2012.10.010.
- 688 Strandén, I., and O.F. Christensen. 2011. Allele coding in genomic evaluation. *Genet. Sel.*
689 *Evol.* 43:25. doi:10.1186/1297-9686-43-25.
- 690 Strandén, I., and D.J. Garrick. 2009. Technical note: Derivation of equivalent computing
691 algorithms for genomic predictions and reliabilities of animal merit. *J. Dairy Sci.*
692 92:2971–2975. doi:10.3168/jds.2008-1929.

Multi-population genomic prediction

- 693 Vandenplas, J., F.G. Colinet, and N. Gengler. 2014. Unified method to integrate and blend
694 several, potentially related, sources of information for genetic evaluation. *Genet. Sel.
695 Evol.* 46:59.
- 696 Vandenplas, J., F.G. Colinet, G. Glorieux, C. Bertozzi, and N. Gengler. 2015. Integration of
697 external estimated breeding values and associated reliabilities using correlations
698 among traits and effects. *J. Dairy Sci.* 98:9044–9050. doi:10.3168/jds.2015-9894.
- 699 Vandenplas, J., and N. Gengler. 2012. Comparison and improvements of different Bayesian
700 procedures to integrate external information into genetic evaluations. *J. Dairy Sci.*
701 95:1513–1526.
- 702 Vandenplas, J., M. Spehar, K. Potocnik, N. Gengler, and G. Gorjanc. 2017. National single-
703 step genomic method that integrates multi-national genomic information. *J. Dairy Sci.*
704 100:465–478. doi:10.3168/jds.2016-11733.
- 705 VanRaden, P.M. 2008. Efficient methods to compute genomic predictions. *J. Dairy Sci.*
706 91:4414–4423. doi:10.3168/jds.2007-0980.
- 707 VanRaden, P.M., and P.G. Sullivan. 2010. International genomic evaluation methods for dairy
708 cattle. *Genet. Sel. Evol.* 42:7. doi:10.1186/1297-9686-42-7.
- 709 VanRaden, P.M., M.E. Tooker, J.R. Wright, C. Sun, and J.L. Hutchison. 2014. Comparison of
710 single-trait to multi-trait national evaluations for yield, health, and fertility. *J. Dairy
711 Sci.* 97:7952–7962.
- 712 Vilhjálmsson, B.J., J. Yang, H.K. Finucane, A. Gusev, S. Lindström, et al. 2015. Modeling
713 linkage disequilibrium increases accuracy of polygenic risk scores. *Am. J. Hum.
714 Genet.* 97:576–592. doi:10.1016/j.ajhg.2015.09.001.

Multi-population genomic prediction

- 715 Wang, H., I. Misztal, I. Aguilar, A. Legarra, and W.M. Muir. 2012. Genome-wide association
716 mapping including phenotypes from relatives without genotypes. *Genet. Res.* 94:73–
717 83. doi:10.1017/S0016672312000274.
- 718 Whittaker, J.C., R. Thompson, and M.C. Denham. 2000. Marker-assisted selection using ridge
719 regression. *Genet. Res.* 75:249–252.
- 720 Wientjes, Y.C.J., P. Bijma, R.F. Veerkamp, and M.P.L. Calus. 2016. An equation to predict
721 the accuracy of genomic values by combining data from multiple traits, populations,
722 or environments. *Genetics* 202:799–823. doi:10.1534/genetics.115.183269.
- 723 Wray, N.R., J. Yang, B.J. Hayes, A.L. Price, M.E. Goddard, et al. 2013. Pitfalls of predicting
724 complex traits from SNPs. *Nat. Rev. Genet.* 14:507–515. doi:10.1038/nrg3457.
- 725 Yang, J., B. Benyamin, B.P. McEvoy, S. Gordon, A.K. Henders, et al. 2010. Common SNPs
726 explain a large proportion of the heritability for human height. *Nat. Genet.* 42:565–
727 569. doi:10.1038/ng.608.
- 728 Yang, J., T. Ferreira, A.P. Morris, S.E. Medland, G.I. of An.T. (GIANT) Consortium, et al.
729 2012. Conditional and joint multiple-SNP analysis of GWAS summary statistics
730 identifies additional variants influencing complex traits. *Nat. Genet.* 44:369–375.
731 doi:10.1038/ng.2213.

732

Multi-population genomic prediction

733 **Table 1** – Genomic prediction accuracy for joint and separate analyses in scenarios with
 734 single or weighted phenotype records and low or high diversity (values are averages across
 735 the five replicates¹)

Phenotypes	Diversity	Analysis	Populations			
			A	B	C	D
Single	Low	Joint	0.811	0.811	0.823	0.815
		Separate	0.705	0.708	0.718	0.718
	High	Joint	0.687	0.686	0.687	0.684
		Separate	0.536	0.537	0.528	0.528
Weighted	Low	Joint	0.860	0.865	0.865	0.862
		Separate	0.720	0.739	0.724	0.727

736 ¹ Standard errors are between 0.003 and 0.016.

737

738

Multi-population genomic prediction

739 **Table 2** - Comparison of estimated allele substitution effects from different analyses with
 740 estimates from the joint statistical analysis using single phenotype records in scenarios with
 741 low and high diversity (values are averages across the five replicates¹)

Analysis	Low diversity		High diversity	
	Correlation	Regression	Correlation	Regression
Separate A	0.71	1.09	0.65	1.10
Separate B	0.71	1.09	0.65	1.10
Separate C	0.71	1.09	0.65	1.11
Separate D	0.71	1.09	0.64	1.10
PEC	1.00	1.00	1.00	1.00
PEC _{within chromosome}	0.99	0.98	0.97	0.95
PEV	0.96	0.80	0.96	0.89
LD _{training}	1.00	1.00	0.98	0.97
LD _{validation}	0.96	0.88	0.93	0.84

742 ¹ Standard errors are between 0.00 and 0.01.

743

Multi-population genomic prediction

744 **Table 3** - Comparison of estimated allele substitution effects from different analyses with
 745 estimates from the joint statistical analysis using weighted phenotype records in the scenario
 746 with low diversity (values are averages across the five replicates with standard errors between
 747 brackets)

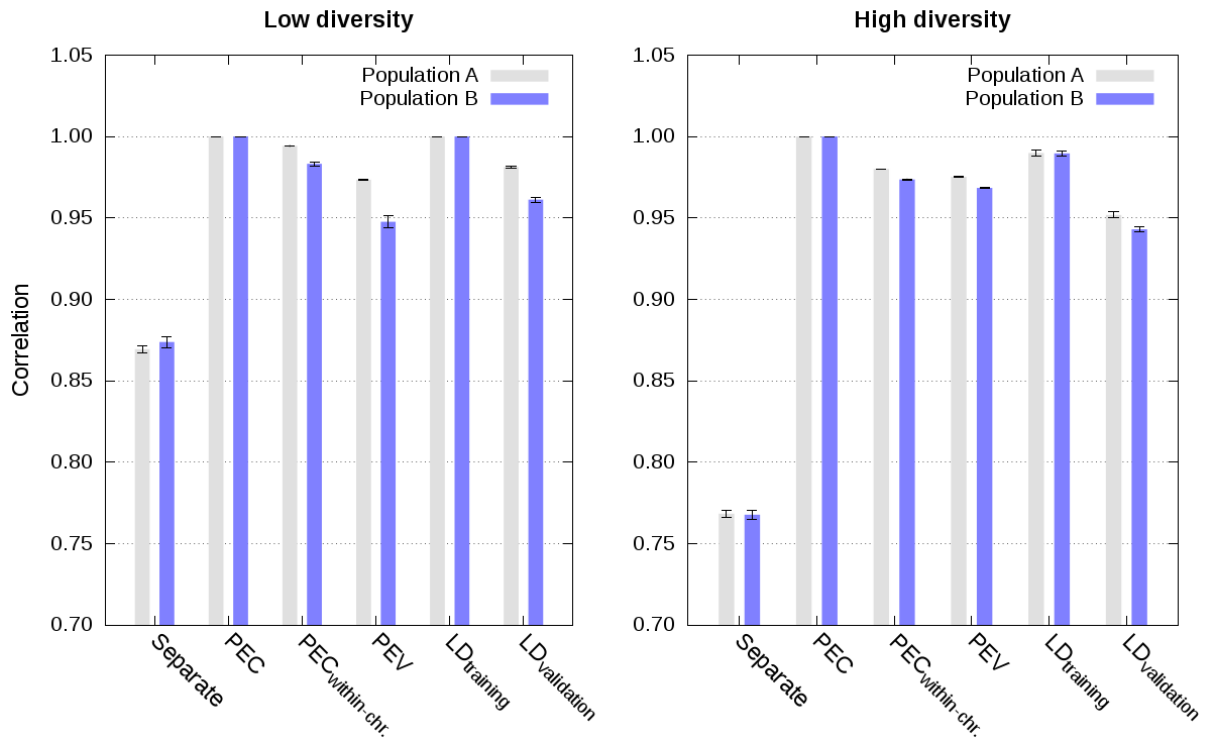
Analysis	Correlation	Regression
Separate A	0.61 (0.10)	0.88 (0.13)
Separate B	0.58 (0.15)	0.62 (0.12)
Separate C	0.56 (0.12)	0.93 (0.23)
Separate D	0.33 (0.08)	0.65 (0.18)
PEC	1.00 (0.00)	0.99 (0.01)
PEC _{within chromosome}	0.96 (0.01)	1.01 (0.02)
PEV	0.92 (0.02)	0.80 (0.05)
LD _{training} (1 marker)	0.77 (0.09)	0.83 (0.10)
LD _{training} (1 chromosome)	0.83 (0.09)	0.95 (0.11)
LD _{validation} (1 marker)	0.73 (0.11)	0.75 (0.13)
LD _{validation} (1 chromosome)	0.71 (0.15)	0.74 (0.18)

748

749

750

FIGURES

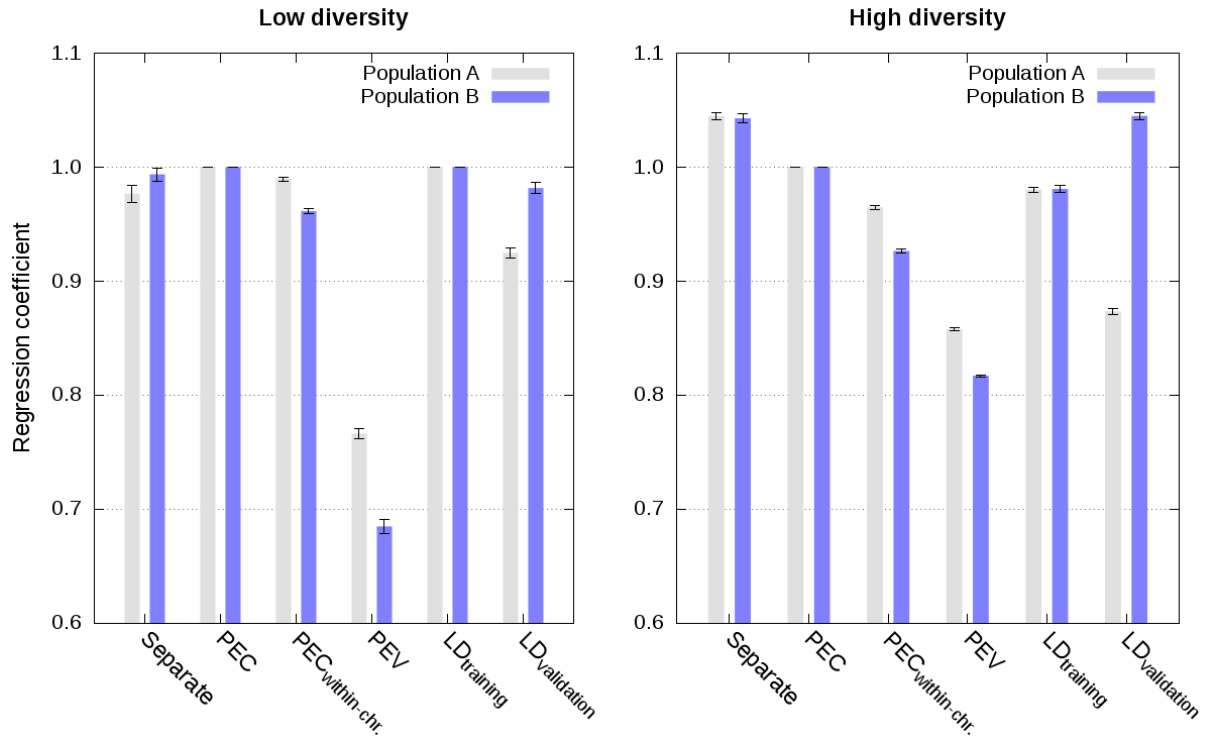


751

752 **Figure 1 - Correlation between estimated genetic values from the joint analysis and from**
 753 **different analyses in populations A and B using a single phenotype record per individual**
 754 **in scenarios with low and high diversity (values are averages across the five replicates**
 755 **with standard errors).**

756

Multi-population genomic prediction



757

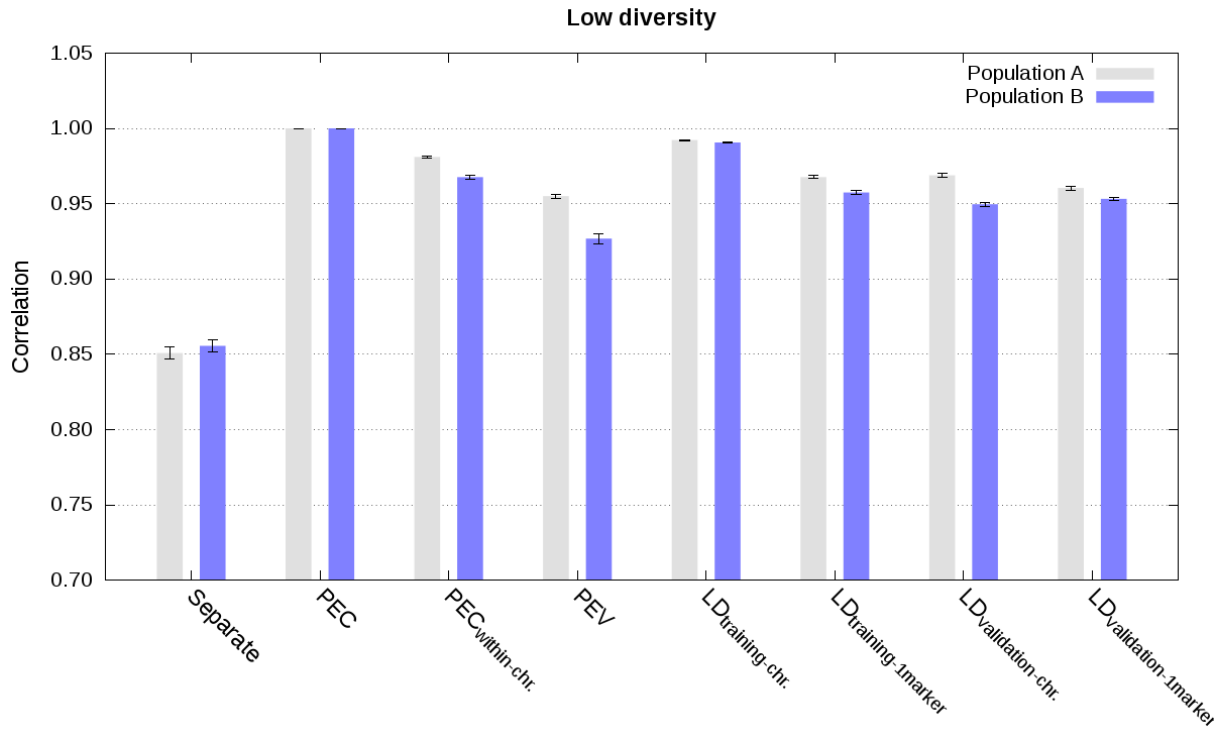
758 **Figure 2 – Regression of estimated genetic values from the joint analysis on estimated**
759 **genetic values from different analyses in populations A and B using a single phenotype**
760 **record per individual in scenarios with low and high diversity (values are averages**
761 **across the five replicates with standard errors).**

762

Multi-population genomic prediction

763

764

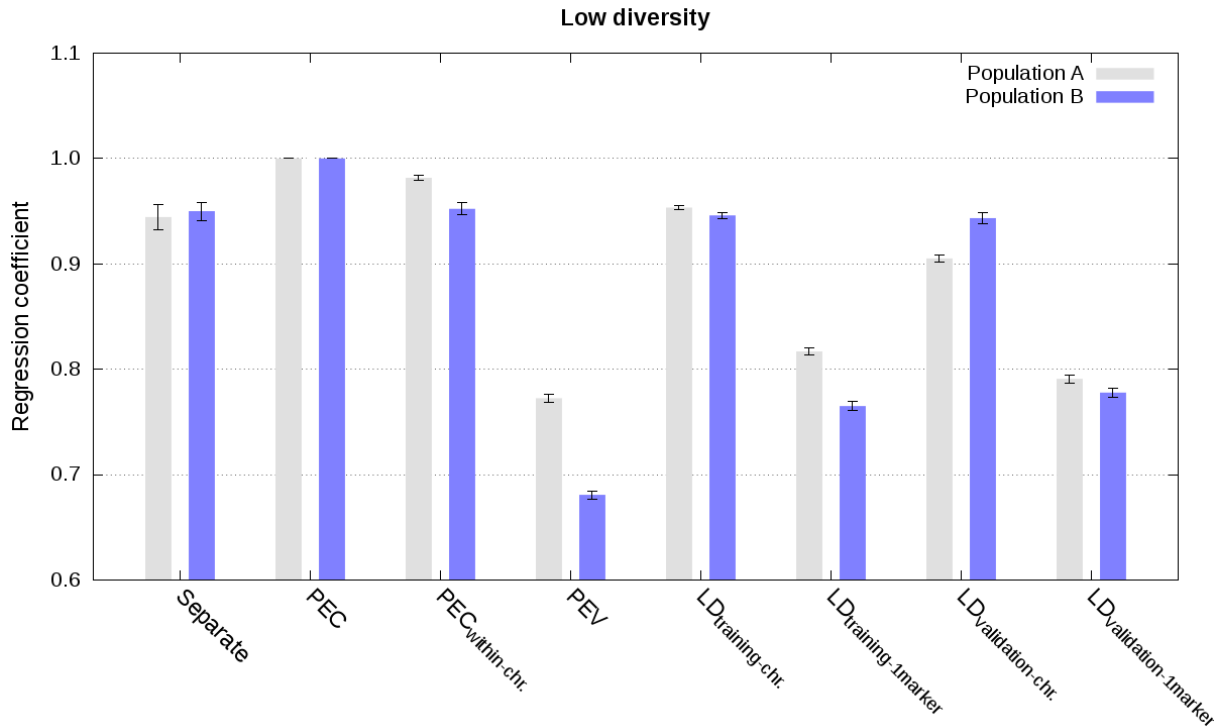


765

766 **Figure 3 - Correlation between estimated genetic values from the joint analysis and from**
767 **different analyses in populations A and B using weighted phenotype records in the**
768 **scenario with low diversity (values are averages across the five replicates with standard**
769 **errors).**

770

Multi-population genomic prediction



771

772 **Figure 4 - Regression of estimated genetic values from the joint analysis on estimated**
773 **genetic values from different analyses in populations A and B using weighted phenotype**
774 **records in the scenario with low diversity (values are averages across the five replicates**
775 **with standard errors).**

776

777 **Appendix A1: Exact integration**

778 Here we detail the derivation of exact integration by means of absorbing the set of
 779 equations that pertain to one dataset. We start with the system of equations for separate analysis
 780 of dataset 1:

$$781 \begin{bmatrix} \mathbf{X}'_1 \mathbf{R}_1^{-1} \sigma_e^{-2} \mathbf{X}_1 & \mathbf{X}'_1 \mathbf{R}_1^{-1} \sigma_e^{-2} \mathbf{Z}_1 \mathbf{W}_1 \\ \mathbf{W}'_1 \mathbf{Z}'_1 \mathbf{R}_1^{-1} \sigma_e^{-2} \mathbf{X}_1 & \mathbf{W}'_1 \mathbf{Z}'_1 \mathbf{R}_1^{-1} \sigma_e^{-2} \mathbf{Z}_1 \mathbf{W}_1 + \mathbf{B}_1^{-1} \sigma_{\alpha_1}^{-2} \end{bmatrix} \begin{bmatrix} \widehat{\boldsymbol{\beta}}_1^* \\ \widehat{\boldsymbol{\alpha}}_1^* \end{bmatrix} = \begin{bmatrix} \mathbf{X}'_1 \mathbf{R}_1^{-1} \sigma_e^{-2} \mathbf{y}_1 \\ \mathbf{W}'_1 \mathbf{Z}'_1 \mathbf{R}_1^{-1} \sigma_e^{-2} \mathbf{y}_1 \end{bmatrix} \quad (\text{A1.1})$$

782 and the system of equations for the joint analysis of datasets 1 and 2:

$$783 \begin{bmatrix} \mathbf{X}'_1 \mathbf{R}_1^{-1} \sigma_e^{-2} \mathbf{X}_1 & \mathbf{0} & \mathbf{X}'_1 \mathbf{R}_1^{-1} \sigma_e^{-2} \mathbf{Z}_1 \mathbf{W}_1 \\ \mathbf{0} & \mathbf{X}'_2 \mathbf{R}_2^{-1} \sigma_e^{-2} \mathbf{X}_2 & \mathbf{X}'_2 \mathbf{R}_2^{-1} \sigma_e^{-2} \mathbf{Z}_2 \mathbf{W}_2 \\ \mathbf{W}'_1 \mathbf{Z}'_1 \mathbf{R}_1^{-1} \sigma_e^{-2} \mathbf{X}_1 & \mathbf{W}'_2 \mathbf{Z}'_2 \mathbf{R}_2^{-1} \sigma_e^{-2} \mathbf{X}_2 & \mathbf{W}'_1 \mathbf{Z}'_1 \mathbf{R}_1^{-1} \sigma_e^{-2} \mathbf{Z}_1 \mathbf{W}_1 + \mathbf{W}'_2 \mathbf{Z}'_2 \mathbf{R}_2^{-1} \sigma_e^{-2} \mathbf{Z}_2 \mathbf{W}_2 + \mathbf{B}_J^{-1} \sigma_{\alpha_J}^{-2} \end{bmatrix} \begin{bmatrix} \widehat{\boldsymbol{\beta}}_1 \\ \widehat{\boldsymbol{\beta}}_2 \\ \widehat{\boldsymbol{\alpha}} \end{bmatrix} =$$

$$784 \begin{bmatrix} \mathbf{X}'_1 \mathbf{R}_1^{-1} \sigma_e^{-2} \mathbf{y}_1 \\ \mathbf{X}'_2 \mathbf{R}_2^{-1} \sigma_e^{-2} \mathbf{y}_2 \\ \mathbf{W}'_1 \mathbf{Z}'_1 \mathbf{R}_1^{-1} \sigma_e^{-2} \mathbf{y}_1 + \mathbf{W}'_2 \mathbf{Z}'_2 \mathbf{R}_2^{-1} \sigma_e^{-2} \mathbf{y}_2 \end{bmatrix}. \quad (\text{A1.2})$$

785 From the first set of equations ($\widehat{\boldsymbol{\beta}}_1$) in (A1.2) it follows:

$$786 \widehat{\boldsymbol{\beta}}_1 = (\mathbf{X}'_1 \mathbf{R}_1^{-1} \sigma_e^{-2} \mathbf{X}_1)^{-1} (\mathbf{X}'_1 \mathbf{R}_1^{-1} \sigma_e^{-2} \mathbf{y}_1 - \mathbf{X}'_1 \mathbf{R}_1^{-1} \sigma_e^{-2} \mathbf{Z}_1 \mathbf{W}_1 \widehat{\boldsymbol{\alpha}}). \quad (\text{A1.3}).$$

787 From the third set of equations ($\widehat{\boldsymbol{\alpha}}$) in (A1.2) it follows:

$$788 \mathbf{W}'_1 \mathbf{Z}'_1 \mathbf{R}_1^{-1} \sigma_e^{-2} \mathbf{X}_1 \widehat{\boldsymbol{\beta}}_1 + \mathbf{W}'_2 \mathbf{Z}'_2 \mathbf{R}_2^{-1} \sigma_e^{-2} \mathbf{X}_2 \widehat{\boldsymbol{\beta}}_2 + (\mathbf{W}'_1 \mathbf{Z}'_1 \mathbf{R}_1^{-1} \sigma_e^{-2} \mathbf{Z}_1 \mathbf{W}_1 +$$

$$789 \mathbf{W}'_2 \mathbf{Z}'_2 \mathbf{R}_2^{-1} \sigma_e^{-2} \mathbf{Z}_2 \mathbf{W}_2 + \mathbf{B}_J^{-1} \sigma_{\alpha_J}^{-2}) \widehat{\boldsymbol{\alpha}} = \mathbf{W}'_1 \mathbf{Z}'_1 \mathbf{R}_1^{-1} \sigma_e^{-2} \mathbf{y}_1 + \mathbf{W}'_2 \mathbf{Z}'_2 \mathbf{R}_2^{-1} \sigma_e^{-2} \mathbf{y}_2. \quad (\text{A1.4}).$$

790 Inserting (A1.3) into (A1.4) gives, after some algebra:

$$791 \mathbf{W}'_2 \mathbf{Z}'_2 \mathbf{R}_2^{-1} \sigma_e^{-2} \mathbf{X}_2 \widehat{\boldsymbol{\beta}}_2 + (\mathbf{W}'_1 \mathbf{Z}'_1 \mathbf{M}_1 \sigma_e^{-2} \mathbf{Z}_1 \mathbf{W}_1 + \mathbf{W}'_2 \mathbf{Z}'_2 \mathbf{R}_2^{-1} \sigma_e^{-2} \mathbf{Z}_2 \mathbf{W}_2 + \mathbf{B}_J^{-1} \sigma_{\alpha_J}^{-2}) \widehat{\boldsymbol{\alpha}}$$

$$792 = \mathbf{W}'_1 \mathbf{Z}'_1 \mathbf{M}_1 \sigma_e^{-2} \mathbf{y}_1 + \mathbf{W}'_2 \mathbf{Z}'_2 \mathbf{R}_2^{-1} \sigma_e^{-2} \mathbf{y}_2$$

793 with $\mathbf{M}_1 = \left(\mathbf{R}_1^{-1} - \mathbf{R}_1^{-1} \mathbf{X}_1 (\mathbf{X}'_1 \mathbf{R}_1^{-1} \mathbf{X}_1)^{-1} \mathbf{X}'_1 \mathbf{R}_1^{-1} \right)$.

794 Now the system of equations (A1.2) can be re-written with the first set of equations

795 $(\widehat{\boldsymbol{\beta}}_1)$ absorbed as:

$$796 \begin{bmatrix} \mathbf{X}'_2 \mathbf{R}_2^{-1} \sigma_e^{-2} \mathbf{X}_2 & \mathbf{X}'_2 \mathbf{R}_2^{-1} \sigma_e^{-2} \mathbf{Z}_2 \mathbf{W}_2 \\ \mathbf{W}'_2 \mathbf{Z}'_2 \mathbf{R}_2^{-1} \sigma_e^{-2} \mathbf{X}_2 & \mathbf{W}'_1 \mathbf{Z}'_1 \mathbf{M}_1 \sigma_e^{-2} \mathbf{Z}_1 \mathbf{W}_1 + \mathbf{W}'_2 \mathbf{Z}'_2 \mathbf{R}_2^{-1} \sigma_e^{-2} \mathbf{Z}_2 \mathbf{W}_2 + \mathbf{B}_J^{-1} \sigma_{\alpha_j}^{-2} \end{bmatrix} \begin{bmatrix} \widehat{\boldsymbol{\beta}}_2 \\ \widehat{\boldsymbol{\alpha}} \end{bmatrix} =$$

$$797 \begin{bmatrix} \mathbf{X}'_2 \mathbf{R}_2^{-1} \sigma_e^{-2} \mathbf{y}_2 \\ \mathbf{W}'_1 \mathbf{Z}'_1 \mathbf{M}_1 \sigma_e^{-2} \mathbf{y}_1 + \mathbf{W}'_2 \mathbf{Z}'_2 \mathbf{R}_2^{-1} \sigma_e^{-2} \mathbf{y}_2 \end{bmatrix}. \quad (\text{A1.4})$$

798 Similarly, the absorption of the first set of equations $(\widehat{\boldsymbol{\beta}}_1^*)$ in separate analysis of dataset

799 1 (A1.1) leads to:

$$800 (\mathbf{W}'_1 \mathbf{Z}'_1 \mathbf{M}_1 \sigma_e^{-2} \mathbf{Z}_1 \mathbf{W}_1 + \mathbf{B}_1^{-1} \sigma_{\alpha_1}^{-2}) \widehat{\boldsymbol{\alpha}}_1^* = \mathbf{W}'_1 \mathbf{Z}'_1 \mathbf{M}_1 \sigma_e^{-2} \mathbf{y}_1, \quad (\text{A1.5})$$

801 where

$$802 \mathbf{W}'_1 \mathbf{Z}'_1 \mathbf{M}_1 \sigma_e^{-2} \mathbf{Z}_1 \mathbf{W}_1 + \mathbf{B}_1^{-1} \sigma_{\alpha_1}^{-2} = \left(PEC(\widehat{\boldsymbol{\alpha}}_1^*) \right)^{-1} \quad (\text{A1.6})$$

803 is the inverse matrix of prediction error covariances of $\widehat{\boldsymbol{\alpha}}_1^*$.

804 Combining (A1.4) and (A1.5) with the use of (A1.6) enables the exact integration of

805 estimates from the separate analysis of dataset 1 into the separate analysis of dataset 2 with the

806 following system of equations:

$$807 \begin{bmatrix} \mathbf{X}'_2 \mathbf{R}_2^{-1} \sigma_e^{-2} \mathbf{X}_2 & \mathbf{X}'_2 \mathbf{R}_2^{-1} \sigma_e^{-2} \mathbf{Z}_2 \mathbf{W}_2 \\ \mathbf{W}'_2 \mathbf{Z}'_2 \mathbf{R}_2^{-1} \sigma_e^{-2} \mathbf{X}_2 & \left(PEC(\widehat{\boldsymbol{\alpha}}_1^*) \right)^{-1} + \mathbf{W}'_2 \mathbf{Z}'_2 \mathbf{R}_2^{-1} \sigma_e^{-2} \mathbf{Z}_2 \mathbf{W}_2 - \mathbf{B}_1^{-1} \sigma_{\alpha_1}^{-2} + \mathbf{B}_J^{-1} \sigma_{\alpha_j}^{-2} \end{bmatrix} \begin{bmatrix} \widehat{\boldsymbol{\beta}}_2 \\ \widehat{\boldsymbol{\alpha}} \end{bmatrix} =$$

$$808 \begin{bmatrix} \mathbf{X}'_2 \mathbf{R}_2^{-1} \sigma_e^{-2} \mathbf{y}_2 \\ \left(PEC(\widehat{\boldsymbol{\alpha}}_1^*) \right)^{-1} \widehat{\boldsymbol{\alpha}}_1^* + \mathbf{W}'_2 \mathbf{Z}'_2 \mathbf{R}_2^{-1} \sigma_e^{-2} \mathbf{y}_2 \end{bmatrix}. \quad (\text{A1.7})$$

809

810 **Appendix A2: Approximate integration**

811 Here we detail the derivation of different approximate integrations by means of
 812 simplified assumptions and use of summary statistics. We start with the expression for
 813 prediction error covariance matrix of allele substitution effects from dataset 1:

$$814 \quad PEC(\widehat{\boldsymbol{\alpha}}_1^*) = (\mathbf{W}'_1 \mathbf{Z}'_1 \mathbf{M}_1 \sigma_e^{-2} \mathbf{Z}_1 \mathbf{W}_1 + \mathbf{B}_1^{-1} \sigma_{\alpha_1}^{-2})^{-1}. \quad (\text{A2.1})$$

815 If we assume that: (1) every individual has a single phenotype record, i.e., $\mathbf{Z}_1 = \mathbf{I}$, (2) residual
 816 variance is homogeneous, i.e. $\mathbf{R}_1 = \mathbf{I}$, and (3) only overall mean is fitted as a fixed effect, i.e.,
 817 $\mathbf{X}_1 = \mathbf{1}$; then we can simplify (A2.1) as:

$$818 \quad PEC(\widehat{\boldsymbol{\alpha}}_1^*) = (\mathbf{W}'_1 \mathbf{Z}'_1 \mathbf{M}_1 \sigma_e^{-2} \mathbf{Z}_1 \mathbf{W}_1 + \mathbf{B}_1^{-1} \sigma_{\alpha_1}^{-2})^{-1},$$

$$819 \quad = (\mathbf{W}'_1 \mathbf{Z}'_1 (\mathbf{R}_1^{-1} - \mathbf{R}_1^{-1} \mathbf{X}_1 (\mathbf{X}'_1 \mathbf{R}_1^{-1} \mathbf{X}_1)^{-1} \mathbf{X}'_1 \mathbf{R}_1^{-1}) \mathbf{Z}_1 \mathbf{W}_1 \sigma_e^{-2} + \mathbf{B}_1^{-1} \sigma_{\alpha_1}^{-2})^{-1},$$

$$820 \quad \approx (\mathbf{W}'_1 (\mathbf{I} - \mathbf{X}_1 (\mathbf{X}'_1 \mathbf{X}_1)^{-1} \mathbf{X}'_1) \mathbf{W}_1 \sigma_e^{-2} + \mathbf{B}_1^{-1} \sigma_{\alpha_1}^{-2})^{-1},$$

$$821 \quad \approx (\mathbf{W}'_1 \mathbf{W}_1 \sigma_e^{-2} + \mathbf{B}_1^{-1} \sigma_{\alpha_1}^{-2})^{-1}, \quad (\text{A2.2})$$

822 because $(\mathbf{I} - \mathbf{X}_1 (\mathbf{X}'_1 \mathbf{X}_1)^{-1} \mathbf{X}'_1) = \mathbf{I} - \mathbf{1}(\mathbf{1}'\mathbf{1})^{-1}\mathbf{1}' = \mathbf{I} - \frac{\mathbf{1}\mathbf{1}'}{n_{ind,1}}$ will tend to the identity matrix
 823 \mathbf{I} with increasing $n_{ind,1}$. The matrix $(\mathbf{I} - \frac{\mathbf{1}\mathbf{1}'}{n_{ind,1}})$, also known as the centering matrix, is a
 824 symmetric and idempotent matrix with off-diagonal elements equal to $-\frac{1}{n_{ind,1}}$ and with
 825 diagonal elements equal to $1 - \frac{1}{n_{ind,1}}$.

826 When genotypes from the dataset 1 are not available, but variance components $\sigma_{\alpha_1}^2$ and
 827 σ_e^2 are, we “only” need to approximate the unknown matrix of genotype sum of squares $\mathbf{W}'_1 \mathbf{W}_1$
 828 in (A2.2). This product can be approximated from linkage-disequilibrium and allele frequency

829 information of the dataset 1, as shown in the following (similarly to Yang et al. (2012),
 830 Vilhjálmsón et al. (2015), and Maier et al. (2018)). Assume that linkage-disequilibrium
 831 between two markers is represented by the correlation of their unphased genotypes (Rogers and
 832 Huff, 2009). Then, a matrix of all pairwise correlations between markers is:

$$833 \quad \mathbf{C} = \left(\text{diag}(\mathbf{T}'_1 \mathbf{T}_1) \right)^{-\frac{1}{2}} \mathbf{T}'_1 \mathbf{T}_1 \left(\text{diag}(\mathbf{T}'_1 \mathbf{T}_1) \right)^{-\frac{1}{2}}, \quad (\text{A2.3})$$

834 where the matrix \mathbf{T}_1 contains centered genotypes of dataset 1 ($\mathbf{T}_1 = \left(\mathbf{I} - \frac{\mathbf{1}\mathbf{1}'}{n_{ind,1}} \right) \mathbf{W}_1 =$
 835 $\mathbf{W}_1 - \frac{1}{n_{ind,1}} \mathbf{1}\mathbf{1}'\mathbf{W}_1$). The matrix product $\mathbf{T}'_1 \mathbf{T}_1$ can be computed as:

$$836 \quad \mathbf{T}'_1 \mathbf{T}_1 = \left(\mathbf{W}_1 - \frac{1}{n_{ind,1}} \mathbf{1}\mathbf{1}'\mathbf{W}_1 \right)' \left(\mathbf{W}_1 - \frac{1}{n_{ind,1}} \mathbf{1}\mathbf{1}'\mathbf{W}_1 \right) = \mathbf{W}'_1 \mathbf{W}_1 - \frac{1}{n_{ind,1}} \mathbf{W}'_1 \mathbf{1}\mathbf{1}'\mathbf{W}_1 -$$

$$837 \quad \frac{1}{n_{ind,1}} \mathbf{W}'_1 \mathbf{1}\mathbf{1}'\mathbf{W}_1 + \frac{1}{n_{ind,1}} \frac{1}{n_{ind,1}} \mathbf{W}'_1 \mathbf{1}\mathbf{1}'\mathbf{1}\mathbf{1}'\mathbf{W}_1 = \mathbf{W}'_1 \mathbf{W}_1 - 4n_{ind,1} \mathbf{p}\mathbf{p}'. \quad (\text{A2.4})$$

838 where $\mathbf{p} = \frac{1}{2n_{ind,1}} \mathbf{W}'_1 \mathbf{1}$ are allele frequencies in dataset 1 (Strandén and Christensen, 2011).

839 Assuming Hardy-Weinberg equilibrium, the i -th diagonal element of the matrix product $\mathbf{T}'_1 \mathbf{T}_1$,
 840 is equivalent to expected genotype sum of squares at the i -th marker, $n_{ind,1} 2p_{i,1}(1 - p_{i,1})$ with
 841 $p_{i,1}$ being the allele frequency of the i -th marker in dataset 1.

842 Combining (A2.3) and (A2.4) we can approximate the unknown matrix of genotype
 843 sum of squares $\mathbf{W}'_1 \mathbf{W}_1$ as:

$$844 \quad \mathbf{W}'_1 \mathbf{W}_1 \approx 4n_{ind,1} \mathbf{p}\mathbf{p}' + \mathbf{V}^{\frac{1}{2}} \mathbf{C} \mathbf{V}^{\frac{1}{2}}, \quad (\text{A2.5})$$

845 where \mathbf{V} is diagonal matrix of expected genotype sum of squares with the i -th diagonal element
 846 equal to $n_{ind,1} 2p_{i,1}(1 - p_{i,1})$.

847

848 **Appendix A3: Estimation of the effective number of records per marker**

849 Here we detail the algorithm for computing the effective number of records per marker
 850 by use of available population parameters (i.e. linkage-disequilibrium, and allele frequency
 851 information) and prediction error variances of $\widehat{\alpha}_1^*$ ($PEV(\widehat{\alpha}_1^*)$) of the dataset 1. We start with the
 852 expression for the prediction error covariance matrix of allele substitution effects from dataset
 853 1:

$$854 \quad PEC(\widehat{\alpha}_1^*) = (\mathbf{W}'_1 \mathbf{Z}'_1 \mathbf{M}_1 \sigma_e^{-2} \mathbf{Z}_1 \mathbf{W}_1 + \mathbf{B}_1^{-1} \sigma_{\alpha_1}^{-2})^{-1}.$$

855 If the number of individuals and the number of records per individual are unknown, we can
 856 assume that a $n_{mar} \times n_{mar}$ diagonal matrix Λ_1 exists such that:

$$857 \quad PEC(\widehat{\alpha}_1^*) \approx \left(\Lambda_1 \left(4\mathbf{p}\mathbf{p}' + \Psi^{\frac{1}{2}} \mathbf{C} \Psi^{\frac{1}{2}} \right) \Lambda_1 \sigma_e^{-2} + \mathbf{B}_1^{-1} \sigma_{\alpha_1}^{-2} \right)^{-1}$$

858 where Ψ is a $n_{mar} \times n_{mar}$ diagonal matrix with the j -th diagonal element equal to
 859 $2p_{j,1}(1 - p_{j,1})$, and the squared j -th diagonal element of Λ_1 represents the effective number of
 860 records for the j -th marker. The term $\left(4\mathbf{p}\mathbf{p}' + \Psi^{\frac{1}{2}} \mathbf{C} \Psi^{\frac{1}{2}} \right)$ is similar to the approximation of the
 861 unknown matrix of genotype sum of squares $\mathbf{W}'_1 \mathbf{W}_1$ (i.e., $\mathbf{W}'_1 \mathbf{W}_1 \approx 4n_{ind,1} \mathbf{p}\mathbf{p}' + \mathbf{V}^{\frac{1}{2}} \mathbf{C} \mathbf{V}^{\frac{1}{2}}$) in
 862 the Appendix A.2. However, it does not involve the number of individuals $n_{ind,1}$ because it is
 863 confounded with the effective number of records.

864 The diagonal matrix Λ_1 can be estimated by solving the nonlinear system of equations

$$865 \quad \text{diag} \left(\left(\Lambda_1 \left(4\mathbf{p}\mathbf{p}' + \Psi^{\frac{1}{2}} \mathbf{C} \Psi^{\frac{1}{2}} \right) \Lambda_1 \sigma_e^{-2} + \mathbf{B}_1^{-1} \sigma_{\alpha_1}^{-2} \right)^{-1} \right) = PEV(\widehat{\alpha}_1^*)$$

866 iteration algorithm (Burden and Faires, 2010) as follows:

$$867 \quad 1) \quad \mathbf{Q}_1^0 = \left(\mathbf{P}^{0^{-1}} - \mathbf{B}_1^{-1} \sigma_{\alpha_1}^{-2} \right) * \left(\text{diag} \left(4\mathbf{p}\mathbf{p}' + \Psi^{\frac{1}{2}} \mathbf{C} \Psi^{\frac{1}{2}} \right) \sigma_e^{-2} \right)^{-1}$$

868 where \mathbf{P}^0 is a diagonal matrix with the i -th diagonal element equal to the PEV of the i -
 869 th marker and $diag\left(4\mathbf{pp}' + \Psi_{\frac{1}{2}}\mathbf{C}\Psi_{\frac{1}{2}}\right)$ contains the diagonal elements of $\left(4\mathbf{pp}' + \Psi_{\frac{1}{2}}\mathbf{C}\Psi_{\frac{1}{2}}\right)$;

871 2) $\Lambda_1^0 = \sqrt{\mathbf{Q}_1^0}$

872 3) $k = 1$

873 4) $\mathbf{P}^k = diag\left(\left(\Lambda_1^{k-1}\left(4\mathbf{pp}' + \Psi_{\frac{1}{2}}\mathbf{C}\Psi_{\frac{1}{2}}\right)\Lambda_1^{k-1}\sigma_e^{-2} + \mathbf{B}_1^{-1}\sigma_{\alpha_1}^{-2}\right)^{-1}\right)$

874 5) $\mathbf{H} = \left(\mathbf{P}^{k-1} - \mathbf{B}_1^{-1}\sigma_{\alpha_1}^{-2}\right) * \left(diag\left(4\mathbf{pp}' + \Psi_{\frac{1}{2}}\mathbf{C}\Psi_{\frac{1}{2}}\right)\sigma_e^{-2}\right)^{-1}$

875 6) $\mathbf{S}^k = \mathbf{Q}_1^0 - \mathbf{H}$

876 7) If trace of \mathbf{S}^k is not sufficiently small:

877 a. $\mathbf{Q}_1^k = \mathbf{Q}_1^{k-1} + \mathbf{H}$

878 b. If any diagonal element in \mathbf{Q}_1^k is negative, set it to 0

879 c. $\Lambda_1^k = \sqrt{\mathbf{Q}_1^k}$

880 d. $k = k + 1$

881 e. Repeat from 4

882 8) $\Lambda_1^k = \sqrt{\mathbf{Q}_1^k}$

883 It is worth noting that the proposed algorithm is similar to algorithms to estimate effective
 884 number of records per individual, where “effective” means that they are free of contributions
 885 from relatives (Miształ and Wiggans, 1988; Vandenplas and Gengler, 2012). The j -th diagonal
 886 element of \mathbf{Q}_1^k can therefore equivalently be considered as the effective number of records for
 887 the j -th marker.

888

889 **Appendix A4: Conversion of allele substitution effects**

890 Here we detail a post-analysis to obtain allele substitution effects estimated using one
 891 type of genotype coding ($\widehat{\alpha}_1^{**}$) by converting estimated genetic values computed for a reference
 892 genotype panel with allele substitution effects for another genotype coding ($\widehat{\alpha}_1^*$). We assume
 893 that allele substitution effects ($\widehat{\alpha}_1^*$) are available with the associated prediction error
 894 (co)variance matrix ($PEC(\widehat{\alpha}_1^*)$), as well as the (co)variance matrix of α_1^* ($Var(\alpha_1^*)$), and
 895 genotypes of a reference panel using a particular type of genotype coding (Γ^*). Estimates of
 896 genetic values for the reference individuals are obtained as $\widehat{\mathbf{g}}_1^* = \Gamma^* \widehat{\alpha}_1^*$.

897 Assuming that estimated genetic values are not influenced by scaling of centered
 898 genotype coding (Strandén and Christensen, 2011; Bouwman et al., 2017), and that the
 899 (co)variances of genetic values are the same irrespective of the genotype coding, we can write
 900 that $\widehat{\mathbf{g}}_1^{**} = \Gamma^{**} \widehat{\alpha}_1^{**} = \widehat{\mathbf{g}}_1^*$ with Γ^{**} being a matrix with reference genotypes using another type
 901 of genotype coding than Γ^* and $\widehat{\mathbf{g}}_1^{**}$ being a vector of estimated genetic values using this type
 902 of genotype coding. Therefore, $\widehat{\alpha}_1^{**}$ can be computed by back-solving as follows (Strandén and
 903 Garrick, 2009; Wang et al., 2012; Bouwman et al., 2017):

$$904 \quad \widehat{\alpha}_1^{**} = \mathbf{B}_1^{**} \Gamma^{**'} (\Gamma^{**} \mathbf{B}_1^{**} \Gamma^{**'})^{-1} \widehat{\mathbf{g}}_1^* = \mathbf{T} \widehat{\mathbf{g}}_1^*$$

905 where \mathbf{B}_1^{**} is a diagonal matrix (e.g., an identity matrix \mathbf{I}) with optional different weights to
 906 differentially shrink different loci.

907 Based on the properties of mixed models (Henderson, 1984), the prediction error
 908 covariance matrix of $\widehat{\alpha}_1^{**}$, $PEC(\widehat{\alpha}_1^{**})$, can be obtained as follows:

Multi-population genomic prediction

$$\begin{aligned}
 909 \quad PEC(\widehat{\boldsymbol{\alpha}}_1^{**}) &= Var(\boldsymbol{\alpha}_1^{**}) - Var(\widehat{\boldsymbol{\alpha}}_1^{**}) = Var(\boldsymbol{\alpha}_1^{**}) - Var(\mathbf{T}\widehat{\mathbf{g}}_1^*) = Var(\boldsymbol{\alpha}_1^{**}) - \mathbf{T}Var(\widehat{\mathbf{g}}_1^*)\mathbf{T}' \\
 910 \quad &= Var(\boldsymbol{\alpha}_1^{**}) - \mathbf{T}\left(Var(\mathbf{g}_1^*) - PEC(\widehat{\mathbf{g}}_1^*)\right)\mathbf{T}' \\
 911 \quad &= Var(\boldsymbol{\alpha}_1^{**}) - \mathbf{T}\left(\boldsymbol{\Gamma}^*Var(\boldsymbol{\alpha}_1^*)\boldsymbol{\Gamma}^{*'} - \boldsymbol{\Gamma}^*PEC(\widehat{\boldsymbol{\alpha}}_1^*)\boldsymbol{\Gamma}^{*'}\right)\mathbf{T}' \\
 912 \quad &= Var(\boldsymbol{\alpha}_1^{**}) - \mathbf{T}\boldsymbol{\Gamma}^*\left(Var(\boldsymbol{\alpha}_1^*) - PEC(\widehat{\boldsymbol{\alpha}}_1^*)\right)\boldsymbol{\Gamma}^{*'}\mathbf{T}'
 \end{aligned}$$

913