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

20    Phone: +31 06 83642304

21

Multi-population genomic prediction

## ABSTRACT

22

23    This study presents a method for genomic prediction that uses individual-level data and

24    summary statistics from multiple populations. Genome-wide markers are nowadays widely

25    used to predict complex traits, and genomic prediction using multi-population data is an

26    appealing approach to achieve higher prediction accuracies. However, sharing of individual-

27    level data across populations is not always possible. We present a method that enables

28    integration of summary statistics from separate analyses with the available individual-level

29    data. The data can either consist of individuals with single or multiple (weighted) phenotype

30    records per individual. We developed a method based on a hypothetical joint analysis model

31    and absorption of population specific information. We show that population specific

32    information is fully captured by estimated allele substitution effects and the accuracy of those

33    estimates, i.e. the summary statistics. The method gives identical result as the joint analysis of

34    all individual-level data when complete summary statistics are available. We provide a series

35    of easy-to-use approximations that can be used when complete summary statistics are not

36    available or impractical to share. Simulations show that approximations enables integration of

37    different sources of information across a wide range of settings yielding accurate predictions.

38    The method can be readily extended to multiple-traits. In summary, the developed method

39    enables integration of genome-wide data in the individual-level or summary statistics form from

40    multiple populations to obtain more accurate estimates of allele substitution effects and

41    genomic predictions.

42

## INTRODUCTION

43

44        Genome-wide markers are nowadays widely used to predict complex traits. This

45  prediction is based on a linear model that partitions for each individual the observed complex

46  phenotype value into systematic effects, comprising at least a population mean, an individual

47  genetic value and an environmental deviation (Fisher, 1918). With genome-wide markers,

48  individual genetic values can be computed from allele substitution effects estimated from

49  individual-level phenotype and genotype data (Meuwissen et al., 2001). Subsequently, genetic

50  values can be also computed for individuals of interest that are genotyped, but not phenotyped.

51  This process is commonly called genomic prediction. In animal and plant breeding, genetic

52  values are used to identify genetically superior individuals and use them as parents of the next

53  generation to improve complex traits like milk yield (Meuwissen et al., 2001; VanRaden, 2008)

54  or grain yield (Schulthess et al., 2016) In human genetics, genetic values can be used to predict

55  individual genetic risk for complex diseases to inform preventive and personalized medicine

56  (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 amount of available data (Daetwyler et al., 2008).

59  To maximize the prediction accuracy, use of all available data is recommended (Henderson,

60  1984; Wray et al., 2013; Vilhjálmsson et al., 2015). In some small populations, collecting large

61  amounts of data is not possible, and a joint analysis across multiple populations is needed to

62  achieve high accuracy (Hozé et al., 2014; Wientjes et al., 2016). However, such joint analysis

63  is often impossible, because of logistic or privacy considerations (Powell and Norman, 1998;

64  Maier et al., 2018). Therefore, several methods were proposed to enable analysis of data from

65  multiple populations when individual-level data is not available (Pasaniuc and Price, 2017; Liu

66  and Goddard, 2018; Maier et al., 2018). These methods approximate a joint analysis by first

67  obtaining summary statistics from separate analyses of individual-level data for each population

4

68  and then combine these summary statistics to estimate genetic values. In human genetics,

69  summary statistics usually consist of publically available allele substitution effects, i.e.,

70  genome-wide associations, together with their standard errors, estimated independently for each

71  marker (Yang et al., 2012; Vilhjálmsson et al., 2015; Maier et al., 2018). In livestock, summary

72  statistics more likely consist of allele substitution effects estimated jointly for all markers,

73  together with prediction error (co)variances (Liu and Goddard, 2018). While these methods

74  may increase prediction accuracy in comparison to separate analyses, a loss in prediction

75  accuracy is expected relative to an analysis using all individual-level data due to approximations

76  (Maier et al., 2018). Further, these methods are based on some assumptions that make them

77  difficult to apply outside their context of development. For example, Maier et al. (2018)

78  implicitly assumed that only a single phenotype record per trait was associated with an

79  individual. While this is usually the case in human genetics, it is not in breeding populations

80  where individuals may have repeated phenotype records for the same trait, e.g., repeated

81  longitudinal production or reproduction records in livestock or replicated field trials in crops,

82  or when phenotype records are measured on a group of individuals and linked to a genotyped

83  relative, e.g., progeny tested bulls for dairy production.

84  The objective of this study was to develop a method that jointly analyses individual-

85  level data and summary statistics from multiple populations with no or limited amount of

86  approximation. The method assumes that individual-level data is composed of marker

87  genotypes and phenotype records that potentially have a variable number of replicates per

88  individual. Further, summary statistics are assumed to be composed of estimated allele

89  substitution effects with an associated measure of accuracy. Different measures of accuracy can

90  be used, which controls the amount of approximation. The developed method is validated with

91  simulated data. The results show that the method enables accurate integration of different

92  sources of information across a wide range of settings.

93

94 **MATERIAL AND METHODS**

95 The first part of this section describes the theory of (1) separate and joint analyses of

96 two individual-level datasets, (2) an exact integration of estimated allele substitution effects

97 from one population into the analysis of another, (3) approximate integrations, and (4)

98 generalization for multiple populations. The second part describes simulations used for

99 validation of the developed method.

100 **Theory**

101 Assume we have two populations with individual-level datasets of phenotyped and

102 genotyped individuals. The two populations and their corresponding datasets are hereafter

103 referred to as 1 and 2. Further assume that both datasets contain the same markers. From this

104 data we want to obtain accurate estimates of allele substitution effects and genetic values for

105 complex traits. We can achieve this by a joint analysis of the two datasets. When one of the

106 datasets is not available, we can achieve this by integrating the results of a separate analysis of

107 the unavailable data into the separate analysis of the available dataset. We show how to perform

108 this integration exactly or approximately.

109 *Separate and joint analyses*

110 A standard marker model, using random regression on marker genotypes, for the

111 separate analysis of dataset $i$ ($i = 1, 2$) is:

112
$$\mathbf{y}_i = \mathbf{X}_i\ \boldsymbol{\beta}_i^* + \mathbf{Z}_i\ \mathbf{W}_i\ \boldsymbol{\alpha}_i^* + \mathbf{e}_i^*, \tag{1}$$

113 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

114 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

115 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 $n_{ind,i} \times$

116 $n_{mar}$ matrix of genotypes $\mathbf{W}_i$, and $\mathbf{e}_i^*$ is the vector $n_{obs,i} \times 1$ of residuals. In this work we

6

117 consider single-nucleotide polymorphism markers, which we code in $\mathbf{W}_i$ as 0 for homozygous

118 aa, 1 for heterozygous aA or Aa, and 2 for homozygous AA. Other genotype coding and

119 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

120 $n_{mar} \times 1$ vector, can be used with no difference in obtained estimates of allele substitution

121 effects (Strandén and Christensen, 2011). We assume a prior multivariate normal (MVN)

122 distribution for allele substitution effects for the separate analyis of the dataset $i$, $\boldsymbol{\alpha}_i^*$, with mean

123 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

124 (e.g., an identity matrix $\mathbf{I}$), and $\sigma_{\alpha_i}^2$ is the variance of allele substitution effects. We also assume

125 that residuals are multivariate normally distributed with mean zero and covariance $\mathbf{R}_i\sigma_e^2$,

126 $\mathbf{e}_i^* \sim MVN(\mathbf{0}, \mathbf{R}_i\ \sigma_e^2)$, where $\mathbf{R}_i$ is a $n_{obs,i} \times n_{obs,i}$ diagonal matrix (e.g., an identity matrix $\mathbf{I}$),

127 and $\sigma_e^2$ is the residual variance. For simplicity and without loss of generality, it is assumed in

128 the following that residual variances are the same for all separate and joint analyses. Variance

129 components $\sigma_{\alpha_i}^2$ and $\sigma_e^2$ are assumed known, as they will have been estimated from the data

130 previously. This marker model is the ridge regression model (Hoerl and Kennard, 1976;

131 Whittaker et al., 2000; Meuwissen et al., 2001; de los Campos et al., 2012) with optional

132 different weights in $\mathbf{B}_i$ (to differentially shrink different loci) and $\mathbf{R}_i$ (to account for

133 heterogeneous residual variance due to variable number of repeated phenotype records per

134 individual).

135 Separate estimates of allele substitution effects $\widehat{\boldsymbol{\alpha}_i^*}$ are obtained by solving the following

136 system of equations:

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

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

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

7

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

141
$$\begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{X}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{X}_2 \end{bmatrix} \begin{bmatrix} \boldsymbol{\beta}_1 \\ \boldsymbol{\beta}_2 \end{bmatrix} + \begin{bmatrix} \mathbf{Z}_1 \ \mathbf{W}_1 \\ \mathbf{Z}_2 \ \mathbf{W}_2 \end{bmatrix} \boldsymbol{\alpha} + \begin{bmatrix} \mathbf{e}_1 \\ \mathbf{e}_2 \end{bmatrix},$$
(3)

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

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

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

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

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

147   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 $n_{obs,i} \times$

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

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

150   system of equations:

151
$$\begin{bmatrix} \mathbf{X}_1'\mathbf{R}_1^{-1}\sigma_e^{-2}\mathbf{X}_1 & \mathbf{0} & \mathbf{X}_1'\mathbf{R}_1^{-1}\sigma_e^{-2}\mathbf{Z}_1 \ \mathbf{W}_1 \\ \mathbf{0} & \mathbf{X}_2'\mathbf{R}_2^{-1}\sigma_e^{-2}\mathbf{X}_2 & \mathbf{X}_2'\mathbf{R}_2^{-1}\sigma_e^{-2}\mathbf{Z}_2 \ \mathbf{W}_2 \\ \mathbf{W}_1'\mathbf{Z}_1'\mathbf{R}_1^{-1}\sigma_e^{-2}\mathbf{X}_1 & \mathbf{W}_2'\mathbf{Z}_2'\mathbf{R}_2^{-1}\sigma_e^{-2}\mathbf{X}_2 & \mathbf{W}_1'\mathbf{Z}_1'\mathbf{R}_1^{-1}\sigma_e^{-2}\mathbf{Z}_1 \ \mathbf{W}_1 \ + \mathbf{W}_2'\mathbf{Z}_2'\mathbf{R}_2^{-1}\sigma_e^{-2}\mathbf{Z}_2 \ \mathbf{W}_2 \ + \mathbf{B}_J^{-1}\sigma_{\alpha_J}^{-2} \end{bmatrix} \begin{bmatrix} \widehat{\boldsymbol{\beta}_1} \\ \widehat{\boldsymbol{\beta}_2} \\ \widehat{\boldsymbol{\alpha}} \end{bmatrix} =$$

152
$$\begin{bmatrix} \mathbf{X}_1'\mathbf{R}_1^{-1}\sigma_e^{-2}\mathbf{y}_1 \\ \mathbf{X}_2'\mathbf{R}_2^{-1}\sigma_e^{-2}\mathbf{y}_2 \\ \mathbf{W}_1'\mathbf{Z}_1'\mathbf{R}_1^{-1}\sigma_e^{-2}\mathbf{y}_1 \ + \mathbf{W}_2'\mathbf{Z}_2'\mathbf{R}_2^{-1}\sigma_e^{-2}\mathbf{y}_2 \end{bmatrix}$$
(4).

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

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

155   ***Exact integration***

156        The integration of estimates of allele substitution effects from one dataset into the

157   analysis of another can be performed by means of absorbing corresponding equations in the

158    joint system of equations. We choose to integrate estimates from the dataset 1 into the analysis

159    of dataset 2. Derivations in Appendix A1 lead to the following system of equations that

160    performs such integration and gives equivalent estimates of allele substitution effects to the

161    joint analysis (4):

162
$$\begin{bmatrix} \mathbf{X}_2'\mathbf{R}_2^{-1}\sigma_e^{-2}\mathbf{X}_2 & \mathbf{X}_2'\mathbf{R}_2^{-1}\sigma_e^{-2}\mathbf{Z}_2\ \mathbf{W}_2 \\ \mathbf{W}_2'\mathbf{Z}_2'\mathbf{R}_2^{-1}\sigma_e^{-2}\mathbf{X}_2 & \left(PEC\big(\widehat{\boldsymbol{\alpha}_1^*}\big)\right)^{-1} + \mathbf{W}_2'\mathbf{Z}_2'\mathbf{R}_2^{-1}\sigma_e^{-2}\mathbf{Z}_2\ \mathbf{W}_2\ -\mathbf{B}_1^{-1}\sigma_{\alpha_1}^{-2} + \mathbf{B}_J^{-1}\sigma_{\alpha_J}^{-2} \end{bmatrix} \begin{bmatrix} \widehat{\boldsymbol{\beta}_2} \\ \widehat{\boldsymbol{\alpha}} \end{bmatrix} =$$

163
$$\begin{bmatrix} \mathbf{X}_2'\mathbf{R}_2^{-1}\sigma_e^{-2}\mathbf{y}_2 \\ \left(PEC\big(\widehat{\boldsymbol{\alpha}_1^*}\big)\right)^{-1}\widehat{\boldsymbol{\alpha}_1^*} + \mathbf{W}_2'\mathbf{Z}_2'\mathbf{R}_2^{-1}\sigma_e^{-2}\mathbf{y}_2 \end{bmatrix},
\tag{5}$$

164    where $\widehat{\boldsymbol{\alpha}_1^*}$ are estimates of allele substitution effects from the separate analysis of dataset 1 using

165    (2), and $\left(PEC\big(\widehat{\boldsymbol{\alpha}_1^*}\big)\right)^{-1}$ is the inverse of the corresponding prediction error covariance (PEC)

166    matrix. The latter can be obtained as $\left(PEC\big(\widehat{\boldsymbol{\alpha}_1^*}\big)\right)^{-1} = \mathbf{W}_1'\mathbf{Z}_1'\mathbf{M}_1\sigma_e^{-2}\mathbf{Z}_1\ \mathbf{W}_1\ + \mathbf{B}_1^{-1}\sigma_{\alpha_1}^{-2}$ with

167    $\mathbf{M}_1 = \left(\mathbf{R}_1^{-1} - \mathbf{R}_1^{-1}\mathbf{X}_1\left(\mathbf{X}_1'\mathbf{R}_1^{-1}\mathbf{X}_1\ \right)^{-1}\mathbf{X}_1'\mathbf{R}_1^{-1}\right)$. Note that only the individual-level dataset 2 and

168    summary statistics from the dataset 1 (i.e., the estimated allele substitution effects and their

169    PEC) are required. Individual-level dataset 1 is therefore not required.

170         It is worth noting that the integration of estimates of allele substitution effects from the

171    dataset 1 into the analysis of dataset 2 can also be obtained from a Bayesian context. Bayes

172    estimators for linear mixed models were discussed by several authors (Lindley and Smith,

173    1972; Dempfle, 1977; Gianola and Fernando, 1986). In a Bayesian context, we can assume

174    the following prior multivariate normal distributions for the marker model (1) applied to

175    dataset 2:

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

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

178 $$[\mathbf{e}_2^* | \mathbf{R}_2 \sigma_e^2] \sim MVN(\mathbf{0}, \mathbf{R}_2 \sigma_e^2).$$

179 Assuming a noninformative prior for $\boldsymbol{\beta}_2^*$, the system of equations (2) for dataset 2 can be

180 obtained by differentiating the joint posterior distribution of $\boldsymbol{\beta}_2^*$ and $\boldsymbol{\alpha}_2^*$ with respect to $\boldsymbol{\beta}_2^*$ and

181 $\boldsymbol{\alpha}_2^*$, and setting the derivatives equal to 0 (Gianola and Fernando, 1986). Integration of

182 estimates of allele substitution effects from dataset 1 into the analysis of dataset 2 can be

183 therefore obained by defining a multivariate normal prior distribution for allele substitution

184 effects in the analysis of dataset 2 using the posterior distribution for allele substitution effects

185 from a separate analysis of dataset 1:

186 $$\left[\boldsymbol{\alpha} | \widehat{\boldsymbol{\alpha}_1^*}, PEC(\widehat{\boldsymbol{\alpha}_1^*}), \mathbf{B}_1 \ \sigma_{\alpha_1}^2, \mathbf{B}_J \ \sigma_{\alpha_J}^2\right] \sim MVN\left(\mathbf{Q}\left(PEC(\widehat{\boldsymbol{\alpha}_1^*})\right)^{-1} \widehat{\boldsymbol{\alpha}_1^*}, \mathbf{Q}\right),$$ (6)

187 $$\mathbf{Q} = \left(\left(PEC(\widehat{\boldsymbol{\alpha}_1^*})\right)^{-1} - \mathbf{B}_1^{-1}\sigma_{\alpha_1}^{-2} + \mathbf{B}_J^{-1}\sigma_{\alpha_J}^{-2}\right)^{-1}.$$

188 The matrix $\mathbf{Q}$ can be considered as the PEC matrix of a hypothetical separate analysis of

189 dataset 1 using the multivariate normal prior distribution for allele substitution effects of the

190 joint analysis, that is $\boldsymbol{\alpha}_1^* \sim MVN\left(\mathbf{0}, \mathbf{B}_J \ \sigma_{\alpha_J}^2\right)$ 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

191 the vector $\mathbf{Q}\left(PEC(\widehat{\boldsymbol{\alpha}_1^*})\right)^{-1} \widehat{\boldsymbol{\alpha}_1^*}$ can be considered as the estimated allele subsitution effects of

192 this hypothectical separate analysis. In animal breeding, a similar approach was used to

193 integrate estimated genetic values and associated accuracies from one genetic evaluation into

194 another genetic evaluation (Quaas and Zhang, 2006; Legarra et al., 2007; Vandenplas and

195 Gengler, 2012).

196 Finally, it is worth noting that the term $\left(PEC(\widehat{\boldsymbol{\alpha}_1^*})\right)^{-1} \widehat{\boldsymbol{\alpha}_1^*}$ can be interpreted as pseudo-

197 phenotypes associated with allele substitution effects of dataset 2, derived from information in

198 dataset 1. In this sense, the system (5) is similar to approaches that compute pseudo-

199    phenotypes from available estimated genetic values where individual-level phenotypic

200    information is not readily available, or is not measured on the individuals themselves but on

201    close relatives. In animal breeding, these approaches are commonly known as deregression of

202    estimated genetic values (Jairath et al., 1998).

203    ***Approximate integration***

204      Exact integration requires the inverse of prediction error covariance matrix from the

205    separate analysis, which could be approximated when unavailable. Genomic analyses of

206    complex traits that combine different datasets commonly have access to estimated allele

207    substitution effects and associated prediction error variances (in different forms), but not the

208    whole prediction error covariance matrix $PEC(\widehat{\boldsymbol{\alpha}_1^*})$ required in (5). We propose several ways

209    to accommodate this situation. We assume that we know, at least, the prediction error variances

210    (PEV) of estimated allele substitution effects $\left(PEV(\widehat{\boldsymbol{\alpha}_1^*})\right)$, the number of individuals $\left(n_{ind,1}\right)$

211    and variance components used in the separate analysis of dataset 1 ($\sigma_{\alpha_1}^2$ and $\sigma_e^2$).

212      When only the prediction error variances of the estimated allele substitution effects

213    $\left(PEV(\widehat{\boldsymbol{\alpha}_1^*})\right)$ are known, while PEC are not, then we can approximate $\left(PEC(\widehat{\boldsymbol{\alpha}_1^*})\right)^{-1}$ with

214    $\left(PEV(\widehat{\boldsymbol{\alpha}_1^*})\right)^{-1}$. This approximation would be accurate if the matrix product $\mathbf{W}_1'\mathbf{W}_1$ has (close

215    to) zero off-diagonal elements, which is dependent on the characteristics of genotypes in dataset

216    1 (e.g., allele frequencies, linkage disequilibrium (LD), and population/family structure). If this

217    is not the case, the approximation will bias the analysis by ignoring off-diagonal elements.

218      When allele frequencies and LD correlations in dataset 1 are known, we can obtain a

219    good approximation of $PEC(\widehat{\boldsymbol{\alpha}_1^*})$ under some conditions (one phenotype record per individual,

220    homogenous residual variance, overall mean is the only fixed effect, and Hardy-Weinberg

221    equilibrium). Derivations in Appendix A2 show that under these conditions we can approximate

11

222 $PEC(\widehat{\boldsymbol{\alpha}_1^*})$ with $\left(\mathbf{W}_1'\mathbf{W}_1 \; \sigma_e^{-2} + \mathbf{B}_1^{-1}\sigma_{\alpha_1}^{-2}\right)^{-1}$ with the unknown matrix $\mathbf{W}_1'\mathbf{W}_1$ approximated

223 from commonly available population parameters (i.e., allele frequencies and LD correlation) as

224 $4n_{ind,1}\mathbf{p}\mathbf{p}' + \mathbf{V}^{\frac{1}{2}}\mathbf{C}\mathbf{V}^{\frac{1}{2}}$, where $\mathbf{p}$ is a $n_{mar} \times 1$ vector of allele frequencies, $\mathbf{V}$ is a $n_{mar} \times n_{mar}$

225 diagonal matrix of expected genotype sum of squares with the $i$-th diagonal element equal to

226 $n_{ind,1}2p_{i,1}(1 - p_{i,1})$, and $\mathbf{C}$ is a $n_{mar} \times n_{mar}$ matrix of pairwise genotype correlations between

227 markers. In practice, the matrix $\mathbf{C}$ for dataset 1 could be unknown, but we can approximate it

228 by using a reference panel that includes, for example, available genotypes of non-phenotyped

229 individuals originating from this population (Yang et al., 2012; Vilhjálmsson et al., 2015; Maier

230 et al., 2018).

231 Finally, we relax the assumption of having a single phenotype record per individual in

232 the preceding approximations. This is relevant when individuals have repeated phenotype

233 records, e.g., repeated longitudinal production or reproduction records in livestock or replicated

234 field trials in crops. A related issue is the violation of assumption of homogenous residual

235 variance when phenotype records are first pre-processed and then used in genomic analyses,

236 e.g., deregressed progeny proofs in livestock (e.g., Garrick et al., 2009) or adjusted field trial

237 means in crops (e.g., Schulz-Streeck et al., 2013; Oakey et al., 2016; Damesa et al., 2017). For

238 these situations, we show in Appendix A3 that we can approximate $PEC(\widehat{\boldsymbol{\alpha}_1^*})$ with

239 $\left(\boldsymbol{\Lambda}_1\left(4\mathbf{p}\mathbf{p}' + \boldsymbol{\Psi}^{\frac{1}{2}}\mathbf{C}\boldsymbol{\Psi}^{\frac{1}{2}}\right)\boldsymbol{\Lambda}_1\sigma_e^{-2} + \mathbf{B}_1^{-1}\sigma_{\alpha_1}^{-2}\right)^{-1}$ where $\boldsymbol{\Psi}$ is a $n_{mar} \times n_{mar}$ diagonal matrix with

240 the $j$-th diagonal element equal to $2p_{j,1}(1 - p_{j,1})$, and $\boldsymbol{\Lambda}_1$ is a $n_{mar} \times n_{mar}$ diagonal matrix

241 with the $j$-th diagonal element representing the square root of effective number of records for

242 the $j$-th marker. The matrix $\boldsymbol{\Lambda}_1$ can be obtained by solving the nonlinear system of equations

243 $diag\left(\left(\boldsymbol{\Lambda}_1\left(4\mathbf{p}\mathbf{p}' + \boldsymbol{\Psi}^{\frac{1}{2}}\mathbf{C}\boldsymbol{\Psi}^{\frac{1}{2}}\right)\boldsymbol{\Lambda}_1\sigma_e^{-2} + \mathbf{B}_1^{-1}\sigma_{\alpha_1}^{-2}\right)^{-1}\right) = PEV(\widehat{\boldsymbol{\alpha}_1^*})$

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

245    It is worth noting that the proposed algorithm requires the inversion of a $n_{mar} \times n_{mar}$ dense

246    matrix at each iteration. This computational cost can be reduced by performing the algorithm

247    for each chromosome separately.

### *Integration with multiple populations*

249    When more than two populations or datasets are available the developed methods can

250    be easily extended. With $n$ datasets, the prior distribution for allele substitution effects in the

251    separate analysis of the $n$-th dataset is defined using the posterior distributions for allele

252    substitution effects from the separate analyses of $n-1$ datasets:

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

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

### **Simulations**

256    We tested developed methods with simulated data that either had low or high genetic

257    diversity. The data was simulated in 5 replicates with the AlphaSim program, which uses the

258    coalescent method for simulation of base population chromosomes and the gene drop method

259    for simulation of chromosome inheritance within a pedigree (Hickey and Gorjanc, 2012; Faux

260    et al., 2016).

261    A diploid genome was simulated with 30 chromosomes, each $10^8$ base pairs long.

262    Coalescent mutation and recombination rate per base pair were set to $10^{-8}$, while effective

263    population size was modelled over time to mimic population history of a livestock population

264    in line with the values reported by MacLeod et al. (2013). Specifically, for the low diversity

265    scenario effective population size of the base population was set to 100 and increased to 120,

13

266 250, 350, 1,000, 1,500, 2,000, 2,500, 3,500, 7,000, 10,000, 17,000, and 62,000 at respectively

267 6, 12, 18, 24, 154, 454, 654, 1,754, 2,354, 3,354, 33,154, and 933,154 generations ago. For the

268 high diversity scenario, effective population size of the base population was set to 10,000 and

269 increased above this value in the same way as in the low diversity scenario; to 17,000 and

270 62,000 at 33,154, and 933,154 generations ago. For each chromosome 10,000 whole

271 chromosome haplotypes were sampled, which on average hosted about 700,000 markers (21

272 million per genome) for the low diversity scenario and 1,400,000 markers (42 million per

273 genome) for the high diversity scenario. Out of these loci 100 per chromosome (3,000 per

274 genome) were sampled as causal loci affecting a complex trait. The allele substitution effect of

275 causal loci was sampled from a normal distribution with mean zero and variance 1/3,000. The

276 effects were used to simulate a complex trait with additive genetic architecture. In addition,

277 2,000 loci per chromosome (60,000 per genome) were selected as markers with the restriction

278 of having minor allele frequency above 0.05.

279       From the base population, founder genomes for four populations (A, B, C, and D) were

280 obtained by random sampling of chromosomes with recombination. The populations were

281 ancestrally related through the common base population, but otherwise maintained

282 independently, i.e., there was no migration between the four populations. Each population was

283 initiated with 10,000 founders (half males and half females) and maintained for 7 generations

284 with constant size. In the low diversity scenario, with the effective population size of 100, 25

285 males and 5,000 females were selected as parents of each generation, while in the high diversity

286 scenario, with the effective population size of 10,000, all 5,000 males and 5,000 females were

287 used. The 25 males were selected on true genetic value, assuming accurate progeny test was

288 available.

289       For every individual in the population we simulated two types of phenotypes. First, an

290 own single phenotype was simulated as the sum of the true genetic value and a residual sampled

291    from a normal distribution with mean zero and residual variance scaled relative to the variance

292    of true genetic value in the base population such that heritability was 0.3. These simulated single

293    phenotype records mimic records measured on the individual. Second, a weighted phenotype

294    was simulated as the sum of the true genetic value and the mean of $n_{weight}$ residuals. Each

295    residual was sampled from a normal distribution with mean zero and residual variance scaled

296    relative to the variance of true genetic value in the base population such that heritability was

297    0.3. The weight $n_{weight}$ was equal to $n_{weight} = 1 + val$ where the real value $val$ was sampled

298    from a geometric distribution with a probability of 0.15. The average $n_{weight}$ was 6.6. These

299    weighted phenotypes mimic either repeated records of an individual or records on multiple

300    progeny of an individual. To satisfy the assumption of identical residual variance across all

301    analyses, phenotype records were divided by the residual standard deviation specific for each

302    population, such that $\sigma_e^2 = 1$. For every individual in each population we stored the true genetic

303    value, own single and weighted phenotype records, associated weight, and 60,000 marker

304    genotypes.

**Analysis**

306        The data was analysed in several ways to evaluate the developed methods. In each case

307    the aim was to obtain accurate genetic values utilizing all the available information.

308    Specifically, we integrated results from separate analysis of populations B, C, and D, into the

309    analysis of population A. We assumed throughout that variance components were known and

310    equal to the rescaled variances. We analysed three scenarios in total. The first and second

311    scenario used population specific training data of randomly sampled 30,000 individuals with

312    single phenotype record from generations 1 to 6 under low and high diversity settings. The third

313    scenario used population specific training data of randomly sampled 10,000 individuals with

314    weighted phenotype record from generations 1 to 6 under low diversity setting. In all scenarios

15

315   all of the 10,000 individuals from generation 7 of each population were considered as validation

316   individuals. The following analyses were performed:

317   1)  A joint analysis of four populations. This was the reference that the other analyses

318       were compared against;

319   2)  A separate analysis for each of the four populations;

320   3)  An exact integration of separate analyses of populations B, C, and D, into the

321       analysis of population A;

322   4)  The same as 3), but approximating the PEC matrix with a partial PEC matrix for

323       each chromosome, i.e., PEC between markers on different chromosomes were set

324       to zero;

325   5)  The same as 3), but approximating the PEC matrix with a diagonal PEV matrix, i.e.,

326       PEC between all markers were set to zero;

327   6)  The same as 3), but approximating the PEC matrix with PEV, allele frequencies,

328       and LD correlations between markers obtained from the training sets. For the

329       scenario with weighted phenotype records, the algorithm for estimating the effective

330       number of records per marker was performed for each marker separately and for

331       each chromsome separately.

332   7)  The same as 6), but with LD correlations between markers computed from

333       validation individuals instead of the training data.

334       For each analysis we calculated genomic prediction accuracy as the Pearson correlation

335   between the true and estimated genetic value in validation individuals. Further, we evaluated

336   the different integrations by comparing estimated genetic values of validation individuals

337   against the estimated genetic values obtained from the joint analysis, which was considered as

338   the reference because it used information from all populations. If integration was fully accurate,

339   there should be no difference between the joint analysis and the analysis with integration. We

16

340    assessed this by (a) accuracy of integration as a Pearson correlation between estimated genetic

341    values from the joint analysis and the analysis with integration (desired value equals 1), (b)

342    calibration of integration as a regression of estimated genetic values from the joint analysis on

343    estimated genetic valuesfrom analysis with integration (desired value equals 1), and (c)

344    magnitude of error in integration as a mean square error (MSE) between estimated genetic

345    values from the joint analysis and from the analysis with integration (desired value equals 0).

346    **Data availability**

347        Supplemental figures are available in File S1. A description of the simulated genotype

348    and phenotype datasets for each scenario is provided in File S2. Simulated genotype and

349    phenotype datasets for the 5 replicates of each scenario are provided in Files S3, S4, and S5.

350    All files were uploaded to Figshare.

351

352          **RESULTS**

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

354          Joint analysis increased genomic prediction accuracy in comparison to separate

355 analyses. This is shown in Table 1. Analysing separately the four datasets gave accuracies of

356 about 0.71 (low diversity) and 0.53 (high diversity) with single phenotype records, and of about

357 0.73 (low diversity) with weighted phenotype records. Analysing jointly the four datasets

358 increased accuracy by 0.09 absolute points with single phenotype records and by 0.12 absolute

359 points with weighted phenotype records.

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

361          For all scenarios the developed method enabled exact integration when complete PEC

362 matrices were used. Integration of estimated allele substitution effects by means of the complete

363 PEC matrix led to the same estimated genetic values as with the joint analysis, as shown by

364 correlation and regression coefficients of 1, and MSE close to 0 (Figures 1-6; Figures S1-S6).

365 For comparison, correlations between estimated genetic values from separate analyses and joint

366 estimated genetic values were about 0.87 (low diversity) and 0.77 (high diversity) with single

367 phenotype records, and 0.85 (low diversity) with weighted phenotype records.

368          Approximate integration by means of partial PEC matrices for each chromosome, that

369 is ignoring PEC between markers on different chromosomes, gave almost as accurate and

370 calibrated estimated genetic values as the exact integration. This is illustrated in Figures 1-6

371 with correlations higher than 0.96, regression coefficients close to 1, and MSE close to 0.

372 Increasing the diversity slightly deteriorated accuracy and calibration of genomic predictions

373 (Figures 1-3; Figures S1-S3).

18

374    Approximate integrations by means of PEV matrices, that is ignoring PEC between all

375    markers, gave quite accurate, but uncalibrated estimated genetic values. This is shown in

376    Figures 1-6 and in Figures S1-S6. Correlations between joint estimated genetic values and

377    estimated genetic values with integration by means of PEV were between 0.95 and 0.98 with

378    single phenotype records and between 0.93 and 0.95 with weighted phenotype records Despite

379    these correlations close to 1, estimated genetic values were uncalibrated, as depicted by

380    regression coefficients below 0.77 for the low diversity scenarios with single and weighted

381    phenotype records, and below 0.86 for the high diversity scenario with single phenotype records

382    (Figures 2, 5, S2, S5).

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

384    When LD information was derived from training data of other populations, approximate

385    integrations by means of PEV, allele frequencies, and LD information, resulted in highly

386    accurate and well calibrated estimated genetic values with single phenotype records. This is

387    shown in Figures 1-3 (Figures S1-S3). Correlation and regression coefficients were equal to 1

388    for the low diversity scenario. Slightly lower values, but still close to 1, were observed for the

389    high diversity scenario. For both low and high diversity scenarios, MSE were close to 0. In

390    contrast, when LD information was derived from validation data of other populations,

391    approximate integrations gave less accurate and well calibrated estimated genetic values. This

392    is shown in Figures 3-6 (Figures S3-S6). For these scenarios, correlations were equal to at least

393    0.94, and regression coefficients varied between 0.87 and 1.05.

394    For the scenario with weighted phenotype records, approximate integrations by means

395    of LD information from training data of other populations resulted in highly accurate and well

396    calibrated estimated genetic values when sets of markers per chromosome were used to estimate

397    the effective number of records for each marker. Correlations between joint estimated genetic

398    values and estimated genetic values with integration were about 0.99 (Figure 4, Figure S4),

399    regression coefficients were about 0.95 (Figure 5, Figure S5), and MSE were close to 0 (Figure

400    6, Figure S6). Using LD information from the validation data of other populations, instead from

401    the training data of other populations, gave slightly less accurate (correlations higher than 0.95),

402    and moderately less calibrated estimated genetic values (regression coefficients between 0.87

403    and 1.04; Figure 4-6; Figures S4-S6). For both cases, estimating the effective numbers of

404    records per marker, instead of for all markers per chromosome simultaneously, reduced

405    accuracy and calibration of estimated genetic values (Figure 4-5; Figures S4-S5).

**Comparison of estimated allele substitution effects**

407        Correlation and regression coefficients between estimated allele substitution effects

408    from the joint analysis and analysis with integration largely followed patterns of the

409    corresponding values for estimated genetic values (Tables 2-3). Correlation and regression

410    coefficients were close to 1 when the integration of estimated allele substitution effects was by

411    means of the complete PEC matrices. Ignoring PEC between markers on different

412    chromosomes, or ignoring PEC between all markers, reduced correlations to between 0.92 and

413    0.99 (Tables 2-3). Using LD information with PEV led to correlations between joint estimates

414    of allele substitution effects and estimates with integration ranging from 0.71 to 0.83 for the

415    scenario with weighted phenotype records (Tables 2-3).

416

20

417 **DISCUSSION**

418      The results show that the developed method enables accurate and well calibrated

419 estimated genetic values for complex traits using both individual-level data and summary

420 statistics. As expected from theory, the analysis of individual-level data and estimated allele

421 substitution effects from other analyses by means of PEC matrices, yielded the same estimates

422 as the joint analysis of all individual-level data. To our knowledge, this is the first time that

423 individual-level data and summary statistics were analysed simultaneously for genomic

424 predictions. As illustrated by simulations, the combined analysis of multiple datasets may

425 increase genomic prediction accuracy over separate analyses of a single dataset. Unfortunately,

426 combining individual-level data from several sources is generally not feasible for several

427 reasons, e.g., political roadblocks, data protections concerns, or data inconsistencies (Powell

428 and Sieber, 1992; Vilhjálmsson et al., 2015; Maier et al., 2018). However, summary statistics,

429 such as estimates of allele substitution effects and associated measures of accuracy (e.g., PEV),

430 are usually available for exchange. The developed method enables increase in genomic

431 prediction accuracy of complex traits by means of jointly analysing the available individual-

432 level data and summary statistics.

433      Accurate integration of estimated allele substitution effects is possible also when the

434 complete PEC matrix is not available. This is important because computing the exact PEC

435 matrix and exchanging it between analyses might be challenging in some cases. For the vast

436 majority of used marker arrays in animal and plant breeding the calculations and data transfers

437 should be doable. For example, most arrays have between 10,000 and 100,000 markers, for

438 which we need between ~1 and ~80 GB of memory to store the PEC matrix and between a

439 minute and a day to invert it on current computers. For a larger number of markers, commonly

440 used in human genetics, the memory requirements and computing time become prohibitive. The

441 results show that in such cases we can still obtain accurate genomic predictions when the

442   integration is done by means of partial PEC matrices for each chromosome. This is expected

443   since high LD between markers mostly occurs within chromosomes. High LD between markers

444   on different chromosomes may especially occur in structured populations and populations

445   under selection (Farnir et al., 2000; Flint-Garcia et al., 2003; Rostoks et al., 2006). Both of these

446   conditions are present in breeding populations. However, the results suggest that LD between

447   chromosomes can be ignored for the purpose of integration for populations with both low and

448   high diversity. The results also show that we can succesfully integrate estimated allele

449   substitution effects when only PEV and allele frequencies from each population are available

450   together with LD information of a reference genotype panel representative of each population.

451   Assuming that such reference genotype panels are available, only estimated allele substitution

452   effects, associated PEV, and allele frequencies need to be exchanged between populations for

453   such analyses. Similar conclusions were drawn from studies combining only summary statistics

454   obtained from genome-wide association studies to perform multi-trait genomic predictions

455   (Maier et al., 2018).

456         Accurate integration of estimated allele substitution effects is possible irrespective of

457   the diversity of the populations and characteristics of genotypes (e.g., allele frequencies, LD).

458   This is obvious, and confirmed by our results, when integration is perfomed by means of

459   complete PEC matrices. When complete PEC matrices are unavailable, accurate integration is

460   possible if the inverses of the PEC matrices can be approximated accurately from available

461   population parameters (i.e. LD and allele frequency information), whatever the level of

462   diversity and characteristics of the populations, as shown by our results or a study combining

463   summary statistics in human genetics (Maier et al., 2018). In our study, the population

464   parameters obtained from the reference panels adequately reflected the characteristics of the

465   training sets. Future studies should be conducted to assess the impact of suboptimal reference

466    panels. Therefore, the developed method is expected to perform well on any type of data, from

467    animal and plant breeding to human genetics, provided accurate information is available.

468    The developed method has some simplifying assumptions that can be readily relaxed.

469    For example, we assumed that the same genotype coding was used in all populations. This

470    assumption can be relaxed when centered genotype coding (i.e., of the form of $(\mathbf{W}_i - \mathbf{1}\mathbf{v}_i')$) is

471    used because variance component estimates, estimates of allele substitution effects and PEC

472    are the same irrespective of the centering of the genotype coding, provided that the model has

473    a fixed general mean, which is considered in the integration (Strandén and Christensen, 2011).

474    Also, centered and scaled (standardised) genotype coding is often used in human genetics,

475    instead of only centered genotype coding (Yang et al., 2010; Speed et al., 2012; Maier et al.,

476    2018). In practice, estimated genetic values are not influenced by scaling of centered genotype

477    coding (Strandén and Christensen, 2011; Bouwman et al., 2017). Therefore, allele substitution

478    effects estimated using one type of genotype scaling could be obtained from a post-analysis by

479    converting estimated genetic values computed for a reference genotype panel into allele

480    substitution effects for another genotype scaling. Converting estimated genetic values into

481    allele substitution effects is often referred to as back-solving of allele substitution effects

482    (Strandén and Garrick, 2009; Strandén and Christensen, 2011; Wang et al., 2012; Bouwman et

483    al., 2017). Prediction error covariances associated with the converted estimated allele

484    subsitution effects could be derived from the (prediction error) covariances of the estimated

485    genetic values (see derivations in Appendix A4).

486    Allele substitution effects estimated from analyses using different different sets of

487    markers or different residual variances, can be used in the integration as well. The assumption

488    that all individuals were genotyped at the same loci could be considered as fullfilled if small

489    differences in the sets of markers are corrected by assuming zero allele substitution effect and

490    zero accuracy for markers not used in an analysis. When large differences between sets of

491    markers are observed, this assumption can be accomodated following two approaches. A first,

492    post-analysis, approach consists of assuming that estimated genetic values are the same for two

493    different sets of markers, allowing the conversion of estimated allele substitution effects from

494    one set of markers to another set of markers (Liu and Goddard, 2018). The conversion can be

495    performed by back-solving estimated allele substitution effects from estimated genetic values,

496    as proposed previously for different genotype codings, or by applying a marker model to the

497    estimated genetic values with the reference set of markers (Liu and Goddard, 2018). A second

498    approach consists of harmonizing genotype data across populations. This approach must be

499    performed before the analyses, and requires therefore coordination between populations.

500    Harmonization of genotype data could be performed by identifying a subset of markers for

501    which all populations are genotyped, or by genotype imputation (e.g., Marchini and Howie,

502    2010). Finally, the assumption that residual variances were the same in all populations, can be

503    relaxed by noting that separate estimates of allele substitution effects $\widehat{\boldsymbol{\alpha}_i^*}$, obtained by the system

504    of equations (2), can be also obtained by the following different formulations:

$$
\begin{aligned}
\widehat{\boldsymbol{\alpha}_i^*} &= \left(\mathbf{W}_i'\mathbf{Z}_i'\mathbf{M}_i\sigma_{e_i}^2\mathbf{Z}_i\mathbf{W}_i + \mathbf{B}_i^{-1}\sigma_{\alpha_i}^{-2}\right)^{-1}\mathbf{W}_i'\mathbf{Z}_i'\mathbf{M}_i\sigma_{e_i}^2\mathbf{y}_i \\
&= \left(\mathbf{W}_i'\mathbf{Z}_i'\mathbf{M}_i\mathbf{Z}_i\mathbf{W}_i + \mathbf{B}_i^{-1}\lambda\right)^{-1}\mathbf{W}_i'\mathbf{Z}_i'\mathbf{M}_i\mathbf{y}_i \\
&= \left(\mathbf{W}_1'\mathbf{Z}_1'\mathbf{M}_1\sigma_{e_f}^{-2}\mathbf{Z}_i\mathbf{W}_i + \mathbf{B}_1^{-1}\lambda\sigma_{e_f}^{-2}\right)^{-1}\mathbf{W}_1'\mathbf{Z}_1'\mathbf{M}_1\sigma_{e_f}^{-2}\mathbf{y}_i
\end{aligned}
$$

505

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

507    For integration of $\widehat{\boldsymbol{\alpha}_i^*}$, $\left(PEC(\widehat{\boldsymbol{\alpha}_i^*})\right)^{-1}$ must be approximated using the residual variance of the

508    focal population ($\sigma_{e_f}^2$) and the effective numbers of records per marker estimated using variance

509    components of the $i$-th analysis. Another way to relax this assumption is to extend our univariate

510    model to a bivariate model, similarly to methods developed to combine different genetic

511    evaluations in animal breeding (Schaeffer, 1994; Vandenplas et al., 2015). In a bivariate model,

512    one trait would represent individual-level data, while the other trait would represent summary

513    statistics. The genetic correlation between the two traits could be estimated based on a subset

514    of individual-level data available for both datasets or based on summary statistics (Bulik-

515    Sullivan et al., 2015). Such an approach would also allow the integegration of summary

516    statistics expressed on a different scale (e.g., different measure units, trait definitions) than the

517    scale of the focal population (Vandenplas et al., 2015).

518        The developed method can be readily generalized to multi-trait models and is therefore

519    a generalization of previous works that were based on several (implicit) assumptions (Liu and

520    Goddard, 2018; Maier et al., 2018). For example, previous works assumed that no individual-

521    level data were available. It was also (implicitly) assumed that only single phenotype records

522    with homogeneous residual variance (Maier et al., 2018), or that the least-squares part of the

523    separate analyses (Liu and Goddard, 2018), were available for integrating estimated allele

524    substitution effects. Both assumptions lead to simple and accurate approximations of PEC

525    matrices as shown in our study. However, we relax all these assumptions, such that our method

526    can jointly analyse individual-level data and summary statistics, with possibly multiple

527    phenotype records per individual.

528

25

## CONCLUSIONS

529

530     We developed a method for genomic prediction that accurately integrates summary

531     statistics obtained from analyses of separate populations into an analysis of individual-level

532     data. The method accommodates use of multiple phenotype (pseudo-)records per individual,

533     and further extensions have been presented to accommodate for differences in residual

534     variances or genotype codings used in the populations. When complete summary statistics

535     information is available the method gives identical genomic predictions as the joint analysis of

536     individual-level data from all populations. When summary statistics information is not

537     complete we can use a series of approximations that give very accurate and well calibrated

538     genomic predictions.

539

540 **ACKNOWLEDGMENTS**

546

547 **LITERATURE CITED**

548 Bouwman, A.C., B.J. Hayes, and M.P.L. Calus. 2017. Estimated allele substitution effects

549 underlying genomic evaluation models depend on the scaling of allele counts. Genet.

550 Sel. Evol. 49. doi:10.1186/s12711-017-0355-9.

551 Bulik-Sullivan, B., H.K. Finucane, V. Anttila, A. Gusev, F.R. Day, P.-R. Loh, ReproGen

552 Consortium, Psychiatric Genomics Consortium, Genetic Consortium for Anorexia

553 Nervosa of the Wellcome Trust Case Control Consortium 3, L. Duncan, J.R.B. Perry,

554 N. Patterson, E.B. Robinson, M.J. Daly, A.L. Price, and B.M. Neale. 2015. An atlas of

555 genetic correlations across human diseases and traits. Nat. Genet. 47:1236–1241.

556 doi:10.1038/ng.3406.

557 Burden, R.L., and J.D. Faires. 2010. Numerical Analysis. 9 edition. Brooks Cole, Boston,

558 MA.

559 Campos, G. de los, D. Gianola, and D.B. Allison. 2010. Predicting genetic predisposition in

560 humans: the promise of whole-genome markers. Nat. Rev. Genet. 11:880–886.

561 doi:10.1038/nrg2898.

562 de los Campos, G., J.M. Hickey, R. Pong-Wong, H.D. Daetwyler, and M.P.L. Calus. 2012.

563 Whole-genome regression and prediction methods applied to plant and animal

564 breeding. Genetics 193:327–345. doi:10.1534/genetics.112.143313.

565 Daetwyler, H.D., B. Villanueva, and J.A. Woolliams. 2008. Accuracy of predicting the

566 genetic risk of disease using a genome-wide approach. PLoS ONE 3.

567    Damesa, T.M., J. Möhring, M. Worku, and H.-P. Piepho. 2017. One step at a time: Stage-wise

568         analysis of a series of experiments. Agron. J. 109:845–857.

569         doi:10.2134/agronj2016.07.0395.

570    Dempfle, L. 1977. Relation entre BLUP (Best Linear Unbiased Prediction) et estimateurs

571         bayésiens. Genet. Sel. Evol. 9:27–32.

572    Farnir, F., W. Coppieters, J.-J. Arranz, P. Berzi, N. Cambisano, B. Grisart, L. Karim, F.

573         Marcq, L. Moreau, M. Mni, C. Nezer, P. Simon, P. Vanmanshoven, D. Wagenaar, and

574         M. Georges. 2000. Extensive genome-wide linkage disequilibrium in cattle. Genome

575         Res. 10:220–227. doi:10.1101/gr.10.2.220.

576    Faux, A.-M., G. Gorjanc, R.C. Gaynor, M. Battagin, S.M. Edwards, D.L. Wilson, S.J. Hearne,

577         S. Gonen, and J.M. Hickey. 2016. AlphaSim: Software for breeding program

578         simulation. Plant Genome 9.

579    Fisher, R.A. 1918. The correlation between relatives on the supposition of Mendelian

580         inheritance. Philos. Trans. R. Soc. Edinb. 52:399–433.

581    Flint-Garcia, S.A., J.M. Thornsberry, and E.S. Buckler. 2003. Structure of linkage

582         disequilibrium in plants. Annu. Rev. Plant Biol. 54:357–374.

583         doi:10.1146/annurev.arplant.54.031902.134907.

584    Garrick, D.J., J.F. Taylor, and R.L. Fernando. 2009. Deregressing estimated breeding values

585         and weighting information for genomic regression analyses. Genet. Sel. Evol. 41:55.

586         doi:10.1186/1297-9686-41-55.

587    Gianola, D., and R.L. Fernando. 1986. Bayesian methods in animal breeding theory. J. Anim.

588         Sci. 63:217–244.

589    Henderson, C.R. 1984. Applications of Linear Models in Animal Breeding. 2nd ed.

590         University of Guelph, Guelph, ON, Canada.

591    Hickey, J.M., and G. Gorjanc. 2012. Simulated data for genomic selection and genome-wide

592         association studies using a combination of coalescent and gene drop methods. G3

593         2:425–427. doi:10.1534/g3.111.001297.

594    Hoerl, A.E., and R.W. Kennard. 1976. Ridge regression iterative estimation of the biasing

595         parameter. Commun. Stat. - Theory Methods 5:77–88.

596         doi:10.1080/03610927608827333.

597    Hozé, C., S. Fritz, F. Phocas, D. Boichard, V. Ducrocq, and P. Croiseau. 2014. Efficiency of

598         multi-breed genomic selection for dairy cattle breeds with different sizes of reference

599         population. J. Dairy Sci. 97:3918–3929. doi:10.3168/jds.2013-7761.

600    Jairath, L., J.C.M. Dekkers, L.R. Schaeffer, Z. Liu, E.B. Burnside, and B. Kolstad. 1998.

601         Genetic evaluation for herd life in Canada. J. Dairy Sci. 81:550–562.

602    Legarra, A., J.K. Bertrand, T. Strabel, R.L. Sapp, J.P. Sanchez, and I. Misztal. 2007. Multi-

603         breed genetic evaluation in a Gelbvieh population. J. Anim. Breed. Genet. 124:286–

604         295.

605    Lindley, D.V., and A.F.M. Smith. 1972. Bayes estimates for the linear model. J. R. Stat. Soc.

606         Ser. B Methodol. 34:1–41.

607    Liu, Z., and M.E. Goddard. 2018. A SNP MACE model for international genomic evaluation:

608         technical challenges and possible solutions. Page 11.393 in Proceedings of the 11th

609         World Congress on Genetics Applied to Livestock Production, Auckland, New

610         Zeland.

611    MacLeod, I.M., D.M. Larkin, H.A. Lewin, B.J. Hayes, and M.E. Goddard. 2013. Inferring

612        demography from runs of homozygosity in whole-genome sequence, with correction

613        for sequence errors. Mol. Biol. Evol. 30:2209–2223.

614    Maier, R.M., Z. Zhu, S.H. Lee, M. Trzaskowski, D.M. Ruderfer, E.A. Stahl, S. Ripke, N.R.

615        Wray, J. Yang, P.M. Visscher, and M.R. Robinson. 2018. Improving genetic

616        prediction by leveraging genetic correlations among human diseases and traits. Nat.

617        Commun. 9:989.

618    Marchini, J., and B. Howie. 2010. Genotype imputation for genome-wide association studies.

619        Nat. Rev. Genet. 11:499–511. doi:10.1038/nrg2796.

620    Meuwissen, T.H.E., B.J. Hayes, and M.E. Goddard. 2001. Prediction of total genetic value

621        using genome-wide dense marker maps. Genetics 157:1819–1829.

622    Misztal, I., and G.R. Wiggans. 1988. Approximation of prediction error variance in large-

623        scale animal models. J. Dairy Sci. 71(Suppl. 2):27–32.

624    Oakey, H., B. Cullis, R. Thompson, J. Comadran, C. Halpin, and R. Waugh. 2016. Genomic

625        selection in multi-environment crop trials. G3 Bethesda Md 6:1313–1326.

626        doi:10.1534/g3.116.027524.

627    Pasaniuc, B., and A.L. Price. 2017. Dissecting the genetics of complex traits using summary

628        association statistics. Nat. Rev. Genet. 18:117–127. doi:10.1038/nrg.2016.142.

629    Powell, R.L., and H.D. Norman. 1998. Use of multinational data to improve national

630        evaluations of Holstein bulls. J. Dairy Sci. 81:2257–2263. doi:10.3168/jds.S0022-

631        0302(98)75805-9.

632    Powell, R.L., and M. Sieber. 1992. Direct and indirect conversion of bull evaluations for yield

633        traits between countries. J. Dairy Sci. 75:1138–1146.

634    Quaas, R.L., and Z. Zhang. 2006. Multiple-breed genetic evaluation in the US beef cattle

635        context: Methodology. Page CD-ROM Comm. 24-12 in Proceedings of the 8th World

636        Congress on Genetics Applied to Livestock Production, Belo Horizonte, Brazil.

637    Rogers, A.R., and C. Huff. 2009. Linkage Disequilibrium Between Loci With Unknown

638        Phase. Genetics 182:839–844. doi:10.1534/genetics.108.093153.

639    Rostoks, N., L. Ramsay, K. MacKenzie, L. Cardle, P.R. Bhat, et al. 2006. Recent history of

640        artificial outcrossing facilitates whole-genome association mapping in elite inbred

641        crop varieties. Proc. Natl. Acad. Sci. U. S. A. 103:18656–18661.

642        doi:10.1073/pnas.0606133103.

643    Schaeffer, L.R. 1994. Multiple-country comparison of dairy sires. J. Dairy Sci. 77:2671–2678.

644    Schulthess, A.W., Y. Wang, T. Miedaner, P. Wilde, J.C. Reif, and Y. Zhao. 2016. Multiple-

645        trait- and selection indices-genomic predictions for grain yield and protein content in

646        rye for feeding purposes. TAG Theor. Appl. Genet. Theor. Angew. Genet. 129:273–

647        287. doi:10.1007/s00122-015-2626-6.

648    Schulz-Streeck, T., J.O. Ogutu, and H.-P. Piepho. 2013. Comparisons of single-stage and two-

649        stage approaches to genomic selection. Theor. Appl. Genet. 126:69–82.

650        doi:10.1007/s00122-012-1960-1.

651    Speed, D., G. Hemani, M.R. Johnson, and D.J. Balding. 2012. Improved heritability

652        estimation from genome-wide SNPs. Am. J. Hum. Genet. 91:1011–1021.

653        doi:10.1016/j.ajhg.2012.10.010.

654    Strandén, I., and O.F. Christensen. 2011. Allele coding in genomic evaluation. Genet. Sel.

655        Evol. 43:25. doi:10.1186/1297-9686-43-25.

656    Strandén, I., and D.J. Garrick. 2009. Technical note: Derivation of equivalent computing

657        algorithms for genomic predictions and reliabilities of animal merit. J. Dairy Sci.

658        92:2971–2975. doi:10.3168/jds.2008-1929.

659    Vandenplas, J., F.G. Colinet, G. Glorieux, C. Bertozzi, and N. Gengler. 2015. Integration of

660        external estimated breeding values and associated reliabilities using correlations

661        among traits and effects. J. Dairy Sci. 98:9044–9050. doi:10.3168/jds.2015-9894.

662    Vandenplas, J., and N. Gengler. 2012. Comparison and improvements of different Bayesian

663        procedures to integrate external information into genetic evaluations. J. Dairy Sci.

664        95:1513–1526.

665    VanRaden, P.M. 2008. Efficient methods to compute genomic predictions. J. Dairy Sci.

666        91:4414–4423. doi:10.3168/jds.2007-0980.

667    Vilhjálmsson, B.J., J. Yang, H.K. Finucane, A. Gusev, S. Lindström, et al. 2015. Modeling

668        linkage disequilibrium increases accuracy of polygenic risk scores. Am. J. Hum.

669        Genet. 97:576–592. doi:10.1016/j.ajhg.2015.09.001.

670    Wang, H., I. Misztal, I. Aguilar, A. Legarra, and W.M. Muir. 2012. Genome-wide association

671        mapping including phenotypes from relatives without genotypes. Genet. Res. 94:73–

672        83. doi:10.1017/S0016672312000274.

673    Whittaker, J.C., R. Thompson, and M.C. Denham. 2000. Marker-assisted selection using ridge

674        regression. Genet. Res. 75:249–252.

675    Wientjes, Y.C.J., P. Bijma, R.F. Veerkamp, and M.P.L. Calus. 2016. An equation to predict

676        the accuracy of genomic values by combining data from multiple traits, populations,

677        or environments. Genetics 202:799–823. doi:10.1534/genetics.115.183269.


678    Wray, N.R., J. Yang, B.J. Hayes, A.L. Price, M.E. Goddard, and P.M. Visscher. 2013. Pitfalls

679        of predicting complex traits from SNPs. Nat. Rev. Genet. 14:507–515.

680        doi:10.1038/nrg3457.


681    Yang, J., B. Benyamin, B.P. McEvoy, S. Gordon, A.K. Henders, et al. 2010. Common SNPs

682        explain a large proportion of the heritability for human height. Nat. Genet. 42:565–

683        569. doi:10.1038/ng.608.


684    Yang, J., T. Ferreira, A.P. Morris, S.E. Medland, G.I. of An.T. (GIANT) Consortium, et al.

685        2012. Conditional and joint multiple-SNP analysis of GWAS summary statistics

686        identifies additional variants influencing complex traits. Nat. Genet. 44:369–375.

687        doi:10.1038/ng.2213.

688


689

690   **Table 1** – Genomic prediction accuracy for joint and separate analyses in scenarios with

691   single or weighted phenotype records and low or high diversity (values are averages across

692   the five replicates[1])

| 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 |

693   [1] Standard errors are between 0.003 and 0.016.

694

695

696 **Table 2** - Comparison of estimated allele substitution effects from different analyses with

697 estimates from the joint statistical analysis using single phenotype records in scenarios with

698 low and high diversity (values are averages across the five replicates[1])

| 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 |

699 [1] Standard errors are between 0.00 and 0.01.

700

36

701 **Table 3** - Comparison of estimated allele substitution effects from different analyses with

702 estimates from the joint statistical analysis using weighted phenotype records in the scenario

703 with low diversity (values are averages across the five replicates with standard errors between

704 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) |

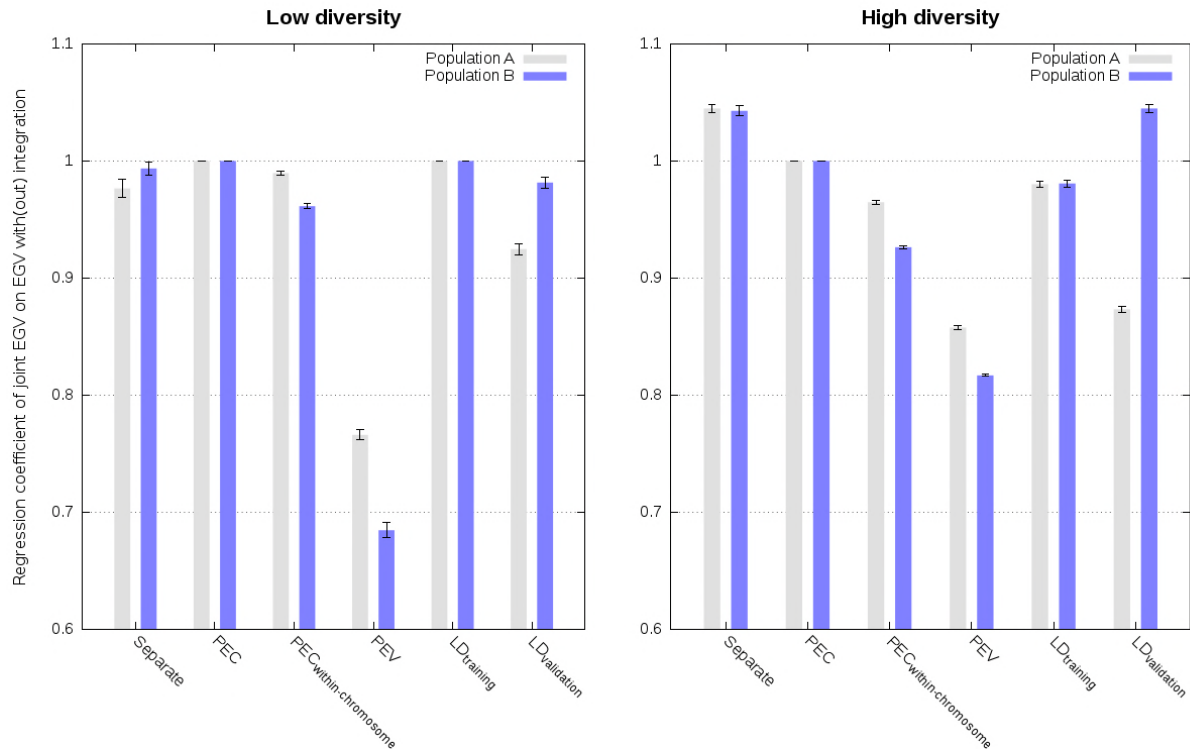705

706

# FIGURES

707



708

709 **Figure 1 - Correlation between estimated genetic values (EGV) from the joint analysis**

710 **and from different analyses in populations A and B using a single phenotype record per**

711 **individual in scenarios with low and high diversity (values are averages across the five**
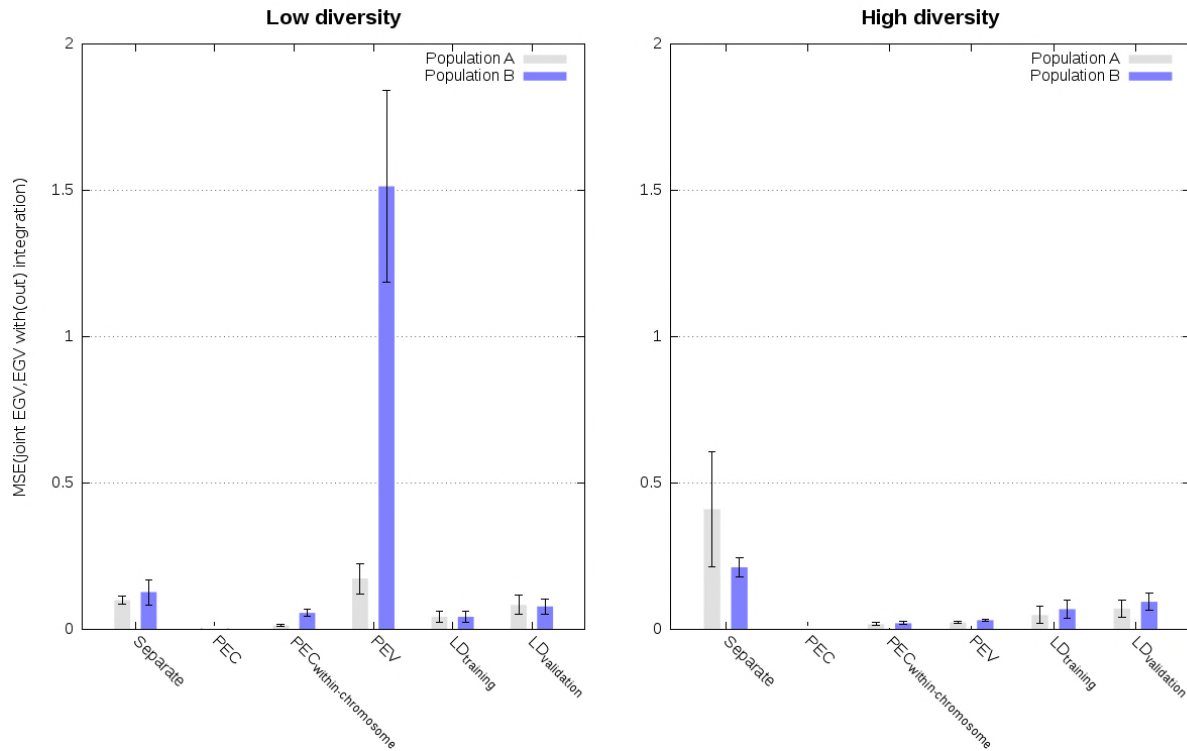
712 **replicates with standard errors).**

713

714

**Figure 2 – Regression of estimated genetic values (EGV) from the joint analysis on estimated genetic values from different analyses in populations A and B using a single phenotype record per individual in scenarios with low and high diversity (values are averages across the five replicates with standard errors).**
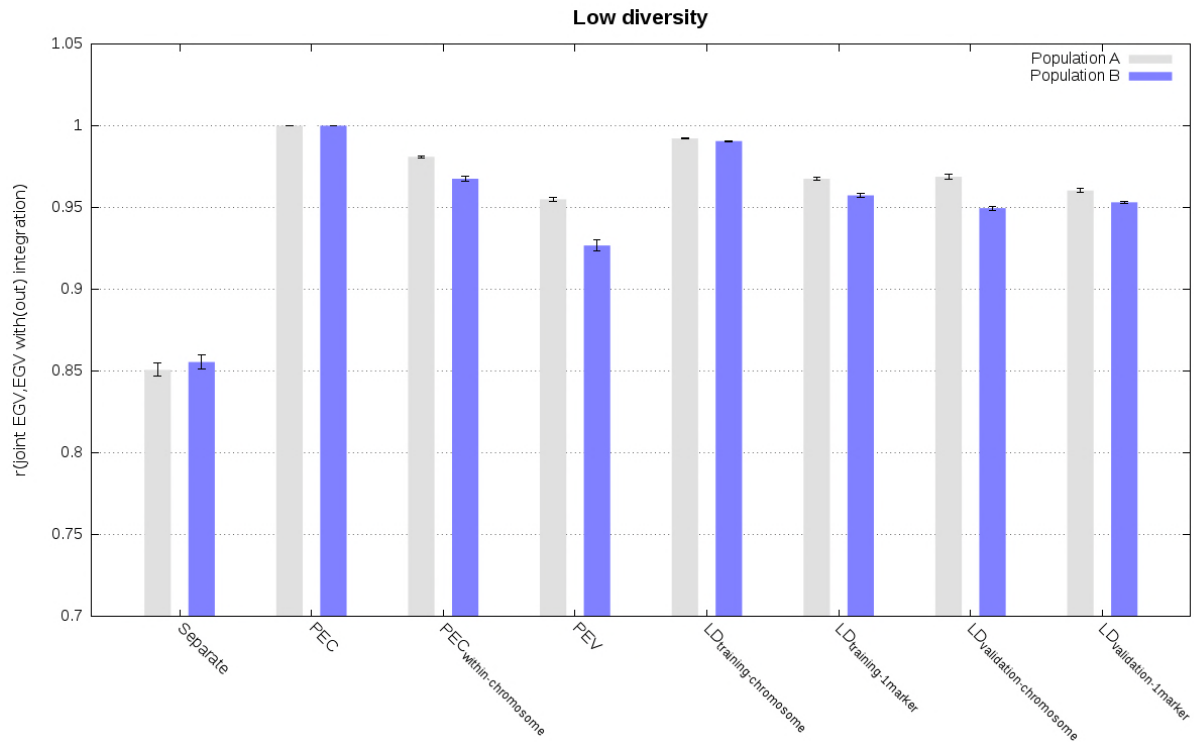
719

**Figure 3 - Mean square errors between joint estimated genetic values (EGV) from the joint analysis and from different analyses in populations A and B using a single phenotype record per individual in scenarios with low and high diversity (values are averages across the five replicates with standard errors).**
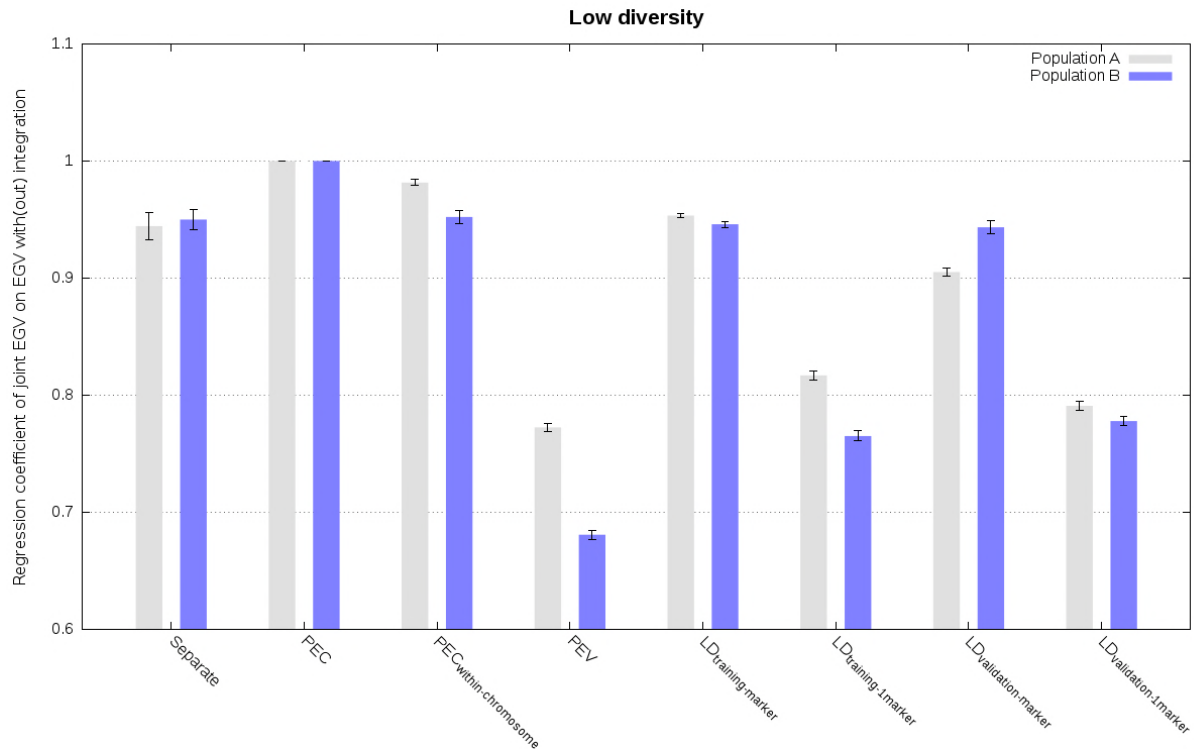
**Figure 4 - Correlation between estimated genetic values (EGV) from the joint analysis and from different analyses in populations A and B using weighted phenotype records in the scenario with low diversity (values are averages across the five replicates with standard errors).**

Multi-population genomic prediction



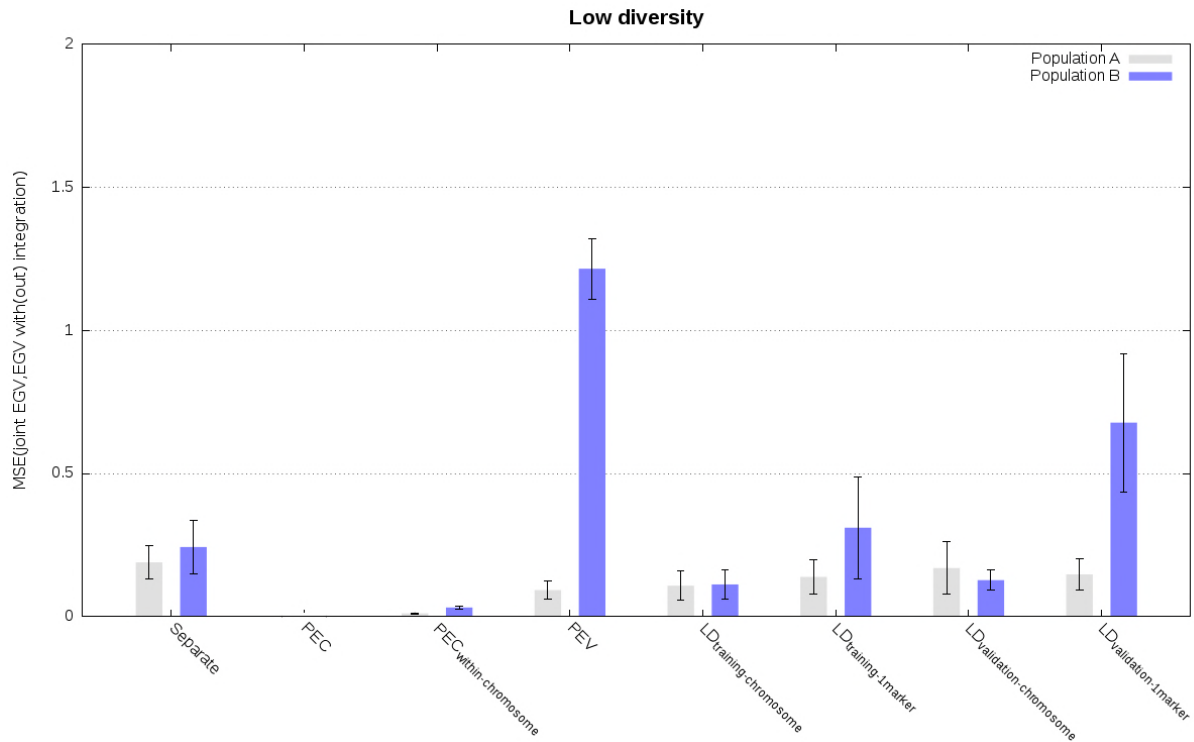**Figure 5 - Regression of estimated genetic values (EGV) from the joint analysis on estimated genetic values from different analyses in populations A and B using weighted phenotype records in the scenario with low diversity (values are averages across the five replicates with standard errors).**

**Figure 6 - Mean square errors (SE) between estimated genetic values (EGV) from the joint analysis and from different analyses in populations A and B using weighted phenotype records in the scenario with low diversity (values are averages across the five replicates with standard errors).**

## Appendix A1: Exact integration

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

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

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

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

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

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

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

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

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

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

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

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

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

44

Multi-population genomic prediction

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

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

765     $\left( \widehat{\boldsymbol{\beta}_1} \right)$ absorbed as:

766
$$\begin{bmatrix} \mathbf{X}_2'\mathbf{R}_2^{-1}\sigma_e^{-2}\mathbf{X}_2 & \mathbf{X}_2'\mathbf{R}_2^{-1}\sigma_e^{-2}\mathbf{Z}_2\,\mathbf{W}_2 \\ \mathbf{W}_2'\mathbf{Z}_2'\mathbf{R}_2^{-1}\sigma_e^{-2}\mathbf{X}_2 & \mathbf{W}_1'\mathbf{Z}_1'\mathbf{M}_1\sigma_e^{-2}\mathbf{Z}_1\,\mathbf{W}_1\ + \mathbf{W}_2'\mathbf{Z}_2'\mathbf{R}_2^{-1}\sigma_e^{-2}\mathbf{Z}_2\,\mathbf{W}_2\ + \mathbf{B}_J^{-1}\sigma_{\alpha_J}^{-2} \end{bmatrix}\begin{bmatrix} \widehat{\boldsymbol{\beta}_2} \\ \widehat{\boldsymbol{\alpha}} \end{bmatrix} =$$

767
$$\begin{bmatrix} \mathbf{X}_2'\mathbf{R}_2^{-1}\sigma_e^{-2}\mathbf{y}_2 \\ \mathbf{W}_1'\mathbf{Z}_1'\mathbf{M}_1\sigma_e^{-2}\mathbf{y}_1\ + \mathbf{W}_2'\mathbf{Z}_2'\mathbf{R}_2^{-1}\sigma_e^{-2}\mathbf{y}_2 \end{bmatrix}. \tag{A1.4}$$

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

769     1 (A1.1) leads to:

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

771     where

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

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

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

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

776     following system of equations:

777
$$\begin{bmatrix} \mathbf{X}_2'\mathbf{R}_2^{-1}\sigma_e^{-2}\mathbf{X}_2 & \mathbf{X}_2'\mathbf{R}_2^{-1}\sigma_e^{-2}\mathbf{Z}_2\,\mathbf{W}_2 \\ \mathbf{W}_2'\mathbf{Z}_2'\mathbf{R}_2^{-1}\sigma_e^{-2}\mathbf{X}_2 & \left( PEC\left( \widehat{\boldsymbol{\alpha}_1^*} \right) \right)^{-1} + \mathbf{W}_2'\mathbf{Z}_2'\mathbf{R}_2^{-1}\sigma_e^{-2}\mathbf{Z}_2\,\mathbf{W}_2\ - \mathbf{B}_1^{-1}\sigma_{\alpha_1}^{-2} + \mathbf{B}_J^{-1}\sigma_{\alpha_J}^{-2} \end{bmatrix}\begin{bmatrix} \widehat{\boldsymbol{\beta}_2} \\ \widehat{\boldsymbol{\alpha}} \end{bmatrix} =$$

778
$$\begin{bmatrix} \mathbf{X}_2'\mathbf{R}_2^{-1}\sigma_e^{-2}\mathbf{y}_2 \\ \left( PEC\left( \widehat{\boldsymbol{\alpha}_1^*} \right) \right)^{-1}\widehat{\boldsymbol{\alpha}_1^*} + \mathbf{W}_2'\mathbf{Z}_2'\mathbf{R}_2^{-1}\sigma_e^{-2}\mathbf{y}_2 \end{bmatrix}. \tag{A1.7}$$

779

**Appendix A2: Approximate integration**

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

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

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

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

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

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

$$\approx \left(\mathbf{W}_1'\mathbf{W}_1\ \sigma_e^{-2} + \mathbf{B}_1^{-1}\sigma_{\alpha_1}^{-2}\right)^{-1}, \tag{A2.2}$$

because $\left(\mathbf{I} - \mathbf{X}_1\left(\mathbf{X}_1'\mathbf{X}_1\right)^{-1}\mathbf{X}_1'\right) = \mathbf{I} - \mathbf{1}(\mathbf{1}'\mathbf{1})^{-1}\mathbf{1}' = \mathbf{I} - \dfrac{\mathbf{11}'}{n_{ind,1}}$ will tend to the identity matrix $\mathbf{I}$ with increasing $n_{ind,1}$. The matrix $\left(\mathbf{I} - \dfrac{\mathbf{11}'}{n_{ind,1}}\right)$, also known as the centering matrix, is a symmetric and idempotent matrix with off-diagonal elements equal to $-\dfrac{1}{n_{ind,1}}$ and with diagonal elements equal to $1 - \dfrac{1}{n_{ind,1}}$.

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

46

799    information of the dataset 1, as shown in the following (similarly to Yang et al. (2012),

800    Vilhjálmsson et al. (2015), and Maier et al. (2018)). Assume that linkage-disequilibrium

801    between two markers is represented by the correlation of their unphased genotypes (Rogers and

802    Huff, 2009). Then, a matrix of all pairwise correlations between markers is:

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

804    where the matrix $\mathbf{T}_1$ contains centered genotypes of dataset 1 ($\mathbf{T}_1 = \left(\mathbf{I} - \frac{\mathbf{1}\mathbf{1}'}{n_{ind,1}}\right)\mathbf{W}_1 =$

805    $\mathbf{W}_1 - \frac{1}{n_{ind,1}}\mathbf{1}\mathbf{1}'\mathbf{W}_1$ ). The matrix product $\mathbf{T}_1'\mathbf{T}_1$ can be computed as:

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

807    $\frac{1}{n_{ind,1}}\mathbf{W}_1'\mathbf{1}\mathbf{1}'\mathbf{W}_1 + \frac{1}{n_{ind,1}}\frac{1}{n_{ind,1}}\mathbf{W}_1'\mathbf{1}\mathbf{1}'\mathbf{1}\mathbf{1}'\mathbf{W}_1 = \mathbf{W}_1'\mathbf{W}_1 - 4n_{ind,1}\mathbf{p}\mathbf{p}'. \tag{A2.4}$

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

809    Assuming Hardy-Weinberg equilibrium, the $i$-th diagonal element of the matrix product $\mathbf{T}_1'\mathbf{T}_1$ ,

810    is equivalent to expected genotype sum of squares at the $i$-th marker, $n_{ind,1}2p_{i,1}\left(1 - p_{i,1}\right)$ with

811    $p_{i,1}$ being the allele frequency of the $i$-th marker in dataset 1.

812        Combining (A2.3) and (A2.4) we can approximate the unknown matrix of genotype

813    sum of squares $\mathbf{W}_1'\mathbf{W}_1$ as:

814    $$\mathbf{W}_1'\mathbf{W}_1 \approx 4n_{ind,1}\mathbf{p}\mathbf{p}' + \mathbf{V}^{\frac{1}{2}}\mathbf{C}\mathbf{V}^{\frac{1}{2}}, \tag{A2.5}$$

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

817

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

819       Here we detail the algorithm for computing the effective number of records per marker

820 by use of available population parameters (i.e. linkage-disequilibrium, and allele frequency

821 information) and prediction error variances of $\widehat{\boldsymbol{\alpha}_1^*}$ ($PEV(\widehat{\boldsymbol{\alpha}_1^*})$) of the dataset 1. We start with the

822 expression for the prediction error covariance matrix of allele substitution effects from dataset

823 1:

824 $PEC(\widehat{\boldsymbol{\alpha}_1^*}) = \left(\mathbf{W}_1'\mathbf{Z}_1'\mathbf{M}_1\sigma_e^{-2}\mathbf{Z}_1\ \mathbf{W}_1\ + \mathbf{B}_1^{-1}\sigma_{\alpha_1}^{-2}\right)^{-1}.$

825 If the number of individuals and the number of records per individual are unknown, we can

826 assume that a $n_{mar} \times n_{mar}$ diagonal matrix $\boldsymbol{\Lambda}_1$ exists such that:

827 $$PEC(\widehat{\boldsymbol{\alpha}_1^*}) \approx \left(\boldsymbol{\Lambda}_1\left(4\mathbf{pp}' + \boldsymbol{\Psi}^{\frac{1}{2}}\mathbf{C}\boldsymbol{\Psi}^{\frac{1}{2}}\right)\boldsymbol{\Lambda}_1\sigma_e^{-2} + \mathbf{B}_1^{-1}\sigma_{\alpha_1}^{-2}\right)^{-1}$$

828 where $\boldsymbol{\Psi}$ is a $n_{mar} \times n_{mar}$ diagonal matrix with the *j*-th diagonal element equal to

829 $2p_{j,1}(1 - p_{j,1})$, and the squared *j*-th diagonal element of $\boldsymbol{\Lambda}_1$ represents the effective number of

830 records for the *j*-th marker. The term $\left(4\mathbf{pp}' + \boldsymbol{\Psi}^{\frac{1}{2}}\mathbf{C}\boldsymbol{\Psi}^{\frac{1}{2}}\right)$ is similar to the approximation of the

831 unknown matrix of genotype sum of squares $\mathbf{W}_1'\mathbf{W}_1$ (i.e., $\mathbf{W}_1'\mathbf{W}_1 \approx 4n_{ind,1}\mathbf{pp}' + \mathbf{V}^{\frac{1}{2}}\mathbf{C}\mathbf{V}^{\frac{1}{2}}$) in

832 the Appendix A.2. However, it does not involve the number of individuals $n_{ind,1}$ because it is

833 confounded with the effective number of records.

834   The diagonal matrix $\boldsymbol{\Lambda}_1$ can be estimated by solving the nonlinear system of equations

835 $diag\left(\left(\boldsymbol{\Lambda}_1\left(4\mathbf{pp}' + \boldsymbol{\Psi}^{\frac{1}{2}}\mathbf{C}\boldsymbol{\Psi}^{\frac{1}{2}}\right)\boldsymbol{\Lambda}_1\sigma_e^{-2} + \mathbf{B}_1^{-1}\sigma_{\alpha_1}^{-2}\right)^{-1}\right) = PEV(\widehat{\boldsymbol{\alpha}_1^*})$ through a fixed-point

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

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

838      where $\mathbf{P}^0$ is a diagonal matrix with the $i$-th diagonal element equal to the PEV of the $i$-

839      th marker and $diag\left(4\mathbf{pp}' + \mathbf{\Psi}^{\frac{1}{2}}\mathbf{C}\mathbf{\Psi}^{\frac{1}{2}}\right)$ contains the diagonal elements of $\left(4\mathbf{pp}' + \right.$

840      $\left.\mathbf{\Psi}^{\frac{1}{2}}\mathbf{C}\mathbf{\Psi}^{\frac{1}{2}}\right)$;

841      2) $\mathbf{\Lambda}_1^0 = \sqrt{\mathbf{Q}_1^0}$

842      3) $k = 1$

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

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

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

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

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

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

849          c. $\mathbf{\Lambda}_1^k = \sqrt{\mathbf{Q}_1^k}$

850          d. $k = k + 1$

851          e. Repeat from 4

852      8) $\mathbf{\Lambda}_1^k = \sqrt{\mathbf{Q}_1^k}$

853 It is worth noting that the proposed algorithm is similar to algorithms to estimate effective

854 number of records per individual, where "effective" means that they are free of contributions

855 from relatives (Misztal and Wiggans, 1988; Vandenplas and Gengler, 2012). The $j$-th diagonal

856 element of $\mathbf{Q}_1^k$ can therefore equivalently be considered as the effective number of records for

857 the $j$-th marker.

858

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

859

860    Here we detail a post-analysis to obtain allele substitution effects estimated using one

861    type of genotype coding ($\widehat{\boldsymbol{\alpha}_1^{**}}$) by converting estimated genetic values computed for a reference

862    genotype panel with allele substitution effects for another genotype coding ($\widehat{\boldsymbol{\alpha}_1^{*}}$). We assume

863    that allele substitution effects ($\widehat{\boldsymbol{\alpha}_1^{*}}$) are available with the associated prediction error

864    (co)variance matrix ($PEC(\widehat{\boldsymbol{\alpha}_1^{*}})$), as well as the (co)variance matrix of $\boldsymbol{\alpha}_1^{*}$ ($Var(\boldsymbol{\alpha}_1^{*})$), and

865    genotypes of a reference panel using a particular type of genotype coding ($\boldsymbol{\Gamma}^{*}$). Estimates of

866    genetic values for the reference individuals are obtained as $\widehat{\mathbf{g}_1^{*}} = \boldsymbol{\Gamma}^{*}\widehat{\boldsymbol{\alpha}_1^{*}}$.

867    Assuming that estimated genetic values are not influenced by scaling of centered

868    genotype coding (Strandén and Christensen, 2011; Bouwman et al., 2017), and that the

869    (co)variances of genetic values are the same irrespective of the genotype coding, we can write

870    that $\widehat{\mathbf{g}_1^{**}} = \boldsymbol{\Gamma}^{**}\widehat{\boldsymbol{\alpha}_1^{**}} = \widehat{\mathbf{g}_1^{*}}$ with $\boldsymbol{\Gamma}^{**}$ being a matrix with reference genotypes using another type

871    of genotype coding than $\boldsymbol{\Gamma}^{*}$ and $\widehat{\mathbf{g}_1^{**}}$ being a vector of estimated genetic values using this type

872    of genotype coding. Therefore, $\widehat{\boldsymbol{\alpha}_1^{**}}$ can be computed by back-solving as follows (Strandén and

873    Garrick, 2009; Wang et al., 2012; Bouwman et al., 2017):

874
$$\widehat{\boldsymbol{\alpha}_1^{**}} = \mathbf{B}_1^{**}\boldsymbol{\Gamma}^{**\prime}(\boldsymbol{\Gamma}^{**}\mathbf{B}_1^{**}\boldsymbol{\Gamma}^{**\prime})^{-1}\widehat{\mathbf{g}_1^{*}} = \mathbf{T}\widehat{\mathbf{g}_1^{*}}$$

875    where $\mathbf{B}_1^{**}$ is a diagonal matrix (e.g., an identity matrix $\mathbf{I}$) with optional different weights to

876    differentially shrink different loci.

877    Based on the properties of mixed models (Henderson, 1984), the prediction error

878    covariance matrix of $\widehat{\boldsymbol{\alpha}_1^{**}}$, $PEC(\widehat{\boldsymbol{\alpha}_1^{**}})$, can be obtained as follows:

879 $$PEC(\widehat{\boldsymbol{\alpha}_1^{**}}) = Var(\boldsymbol{\alpha}_1^{**}) - Var(\widehat{\boldsymbol{\alpha}_1^{**}}) = Var(\boldsymbol{\alpha}_1^{**}) - Var(\mathbf{T}\widehat{\mathbf{g}_1^{*}}) = Var(\boldsymbol{\alpha}_1^{**}) - \mathbf{T}Var(\widehat{\mathbf{g}_1^{*}})\mathbf{T}'$$

880 $$= Var(\boldsymbol{\alpha}_1^{**}) - \mathbf{T}\left(Var(\mathbf{g}_1^{*}) - PEC(\widehat{\mathbf{g}_1^{*}})\right)\mathbf{T}'$$

881 $$= Var(\boldsymbol{\alpha}_1^{**}) - \mathbf{T}(\boldsymbol{\Gamma}^{*}Var(\boldsymbol{\alpha}_1^{*})\boldsymbol{\Gamma}^{*\prime} - \boldsymbol{\Gamma}^{*}PEC(\widehat{\boldsymbol{\alpha}_1^{*}})\boldsymbol{\Gamma}^{*\prime})\mathbf{T}'$$

882 $$= Var(\boldsymbol{\alpha}_1^{**}) - \mathbf{T}\boldsymbol{\Gamma}^{*}\left(Var(\boldsymbol{\alpha}_1^{*}) - PEC(\widehat{\boldsymbol{\alpha}_1^{*}})\right)\boldsymbol{\Gamma}^{*\prime}\mathbf{T}'$$

883