

1 **Properties of genomic relationships for estimating current genetic**
2 **variances within and genetic correlations between populations**

3

4 Yvonne C.J. Wientjes*, Piter Bijma*, J r mie Vandenplas*, Mario P.L. Calus*

5

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

7 Wageningen, The Netherlands

8

9

10

11

12

13

14

15 Running title: Genetic correlation between populations

16

17 Key words: genetic correlation between populations, genomic relationships, genetic variance,

18 multi-trait model

19

20 Author information:

21 Yvonne Wientjes

22 Wageningen University & Research

23 Animal Breeding and Genomics

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

25 E-mail: yvonne.wientjes@wur.nl

26 Phone: +31 317 481 904

27

28

29

30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49

ABSTRACT

Different methods are available to calculate multi-population genomic relationship matrices. Since those matrices differ in base population, it is anticipated that the method used to calculate the genomic relationship matrix affect the estimate of genetic variances, covariances and correlations. The aim of this paper is to define a multi-population genomic relationship matrix to estimate current genetic variances within and genetic correlations between populations. The genomic relationship matrix containing two populations consists of four blocks, one block for population 1, one block for population 2, and two blocks for relationships between the populations. It is known, based on literature, that current genetic variances are estimated when the current population is used as base population of the relationship matrix. In this paper, we theoretically derived the properties of the genomic relationship matrix to estimate genetic correlations and validated it using simulations. When the scaling factors of the genomic relationship matrix fulfill the property $k_{12} = \sqrt{k_1} \sqrt{k_2}$, the genetic correlation is estimated even though estimated variance components are not necessarily related to the current population. When this property is not met, the correlation based on estimated variance components should be multiplied by $\frac{\sqrt{k_1} \sqrt{k_2}}{k_{12}}$ to rescale the genetic correlation. In this study we present a genomic relationship matrix which directly results in current genetic variances as well as genetic correlations between populations.

50

INTRODUCTION

51 When estimating additive genetic values of individuals, the relationships between
52 individuals are used to describe the covariance between additive genetic values for a specific
53 trait. Those covariances between individuals are best represented by the relationships at causal
54 loci. Since causal loci are generally unknown, different approaches have been developed to
55 estimate relationships from genomic marker data (e.g., VanRaden 2008; Powell *et al.* 2010;
56 Yang *et al.* 2010). As long as causal loci and genomic markers have the same properties, such
57 as allele frequency distribution, relationships at the markers are observed to be an unbiased
58 estimate of relationships at the causal loci (Yang *et al.* 2010; Yang *et al.* 2015).

59 Relationships are expressed relative to a base population, consisting of unrelated
60 individuals that have average self-relationships of one, for which the additive genetic variance
61 is estimated. The base population of a genomic relationship matrix depends on the method
62 used to calculate the relationship matrix, therefore estimated variances differ across methods
63 (Speed and Balding 2015; Legarra 2016). By using the current allele frequencies to calculate
64 the genomic relationship matrix, the current population is the base population for which
65 additive genetic variances are estimated (Hayes *et al.* 2009).

66 Genomic relationships can also be calculated between distantly related individuals, for
67 example between individuals from different populations. Those relationships can be used to
68 estimate genetic correlations between populations using a multi-trait model (Karoui *et al.*
69 2012), where the same trait in each population is modelled as a different trait. Due to
70 differences in environments and allele frequencies, in combination with non-additive effects,
71 the allele substitution effects of causal loci can differ between populations (e.g., Fisher 1918;
72 Fisher 1930; Falconer 1952). Moreover, some causal loci might only segregate in one of the
73 populations. Therefore, the genetic correlation between populations can differ from 1.

74 The genetic correlation between populations is an important parameter, since it is used to
75 understand the genetic architecture and evolution of complex traits, such as disease traits in
76 humans (De Candia *et al.* 2013; Brown *et al.* 2016). Moreover, the genetic correlation
77 determines whether information can be shared across populations as done in multi-population
78 genomic prediction (Wientjes *et al.* 2015; Wientjes *et al.* 2016), which is of importance for
79 animals (e.g., Karoui *et al.* 2012; Olson *et al.* 2012), plants (e.g., Lehermeier *et al.* 2015) and
80 humans (e.g., De Candia *et al.* 2013).

81 Different methods are available to calculate multi-population genomic relationship
82 matrices (Harris and Johnson 2010; Erbe *et al.* 2012; Chen *et al.* 2013; Makgahlela *et al.*
83 2013). The two most important differences between the methods are: 1) the assumed relation
84 between effect size and allele frequency of markers; namely assuming effect size and allele
85 frequency are independent (e.g., method 1 of VanRaden (2008)) or assuming that markers
86 with a lower allele frequency have a larger effect (e.g., method 2 of VanRaden (2008) and
87 Yang (2010)), and 2) the allele frequency that is used; namely allele frequencies specific to
88 each population, the average allele frequency across the populations, or the estimated allele
89 frequency when the populations separated. Since relationships between individuals differ
90 across those methods, it is anticipated that the method used to calculate the genomic
91 relationship matrix affects the estimate of the genetic correlation.

92 Therefore, the aim of this paper is to define a multi-population genomic relationship matrix
93 to estimate current genetic variances within and genetic correlations between populations. We
94 theoretically derive a relationship matrix with this property and validate it with simulations.
95 To rule out the effect of differences in linkage disequilibrium between markers and causal
96 loci, we will focus in the entire paper on a situation where causal loci are used to calculate the
97 relationships.

98

MATERIALS AND METHODS

99 Theory

100 The additive genetic correlation, r_g , is the correlation between additive genetic values (A)
 101 for two traits of the same individual (Bohren *et al.* 1966; Falconer and Mackay 1996). In an
 102 additive model and under the assumptions that the correlation originates from pleiotropy,
 103 genetic values are independent between loci, and allele substitution effects are independent
 104 from allele frequency, r_g is equal to the average correlation between allele substitution effects
 105 of the two traits, denoted as trait 1 and 2, at causal loci. This equality can be shown for
 106 individual i by considering both genotypes (z) and allele substitution effects (α) at all n_c
 107 causal loci as random:

$$\begin{aligned}
 108 \quad \text{Var}(A_{i1}) &= \text{Var}\left(\sum_j z_{ij}\alpha_{1j}\right) = E\left(\left(\sum_j z_{ij}\alpha_{1j}\right)\left(\sum_l z_{il}\alpha_{1l}\right)\right) = \sum_j E(z_{ij}z_{ij})E(\alpha_{1j}\alpha_{1j}) \\
 &= n_c E(z_{ij}z_{ij})E(\alpha_{1j}\alpha_{1j}) \\
 109 \quad \text{Var}(A_{i2}) &= n_c E(z_{ij}z_{ij})E(\alpha_{2j}\alpha_{2j}) \\
 110 \quad \text{Cov}(A_{i1}, A_{i2}) &= \text{Cov}\left(\sum_j z_{ij}\alpha_{1j}, \sum_l z_{il}\alpha_{2l}\right) = E\left(\left(\sum_j z_{ij}\alpha_{1j}\right)\left(\sum_l z_{il}\alpha_{2l}\right)\right) = \sum_j E(z_{ij}z_{ij})E(\alpha_{1j}\alpha_{2j}) \\
 &= n_c E(z_{ij}z_{ij})E(\alpha_{1j}\alpha_{2j}) \\
 r_g &= \frac{\text{Cov}(A_{i1}, A_{i2})}{\sqrt{\text{Var}(A_{i1})} \sqrt{\text{Var}(A_{i2})}} = \frac{n_c E(z_{ij}z_{ij})E(\alpha_{1j}\alpha_{2j})}{\sqrt{n_c E(z_{ij}z_{ij}) E(\alpha_{1j}\alpha_{1j})} \sqrt{n_c E(z_{ij}z_{ij}) E(\alpha_{2j}\alpha_{2j})}} \\
 111 \quad &= \frac{E(\alpha_{1j}\alpha_{2j})}{\sqrt{E(\alpha_{1j}\alpha_{1j})} \sqrt{E(\alpha_{2j}\alpha_{2j})}} = \frac{\sigma_{\alpha_{12}}}{\sqrt{\sigma_{\alpha_1}^2} \sqrt{\sigma_{\alpha_2}^2}} = r_\alpha \quad (1)
 \end{aligned}$$

112 where j and l denote the different causal loci. Genotypes are represented by allele counts
 113 coded as 0, 1 and 2 that are centered by subtracting $2p$, where p is the allele frequency for the
 114 counted allele.

115 Similar to genetic correlations between traits in one population, the genetic correlation (r_g)
 116 between populations can be estimated in a multi-trait model using a relationship matrix and

117 REML by modelling the phenotypes of two populations as different traits (Karoui *et al.*
 118 2012). This approach is also known as multi-trait GREML. In the following, we will refer to
 119 trait 1 as the trait expressed in population 1 and to trait 2 as the trait expressed in population 2.
 120 When considering performance in different populations as different traits, individuals have a
 121 phenotype for only one trait. Therefore, the (co)variance structure of the additive genetic
 122 values can be written as (Visscher *et al.* 2014):

$$123 \quad \begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \end{bmatrix} \sim N \left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \text{Var}(\mathbf{a}_1) & \text{Cov}(\mathbf{a}_1, \mathbf{a}_2) \\ \text{Cov}(\mathbf{a}_2, \mathbf{a}_1) & \text{Var}(\mathbf{a}_2) \end{bmatrix} \right) = N \left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \mathbf{G}_{11}\sigma_1^2 & \mathbf{G}_{12}\sigma_{12} \\ \mathbf{G}_{21}\sigma_{12} & \mathbf{G}_{22}\sigma_2^2 \end{bmatrix} \right). \quad (2)$$

124 where \mathbf{a}_1 is the vector with additive genetic values for individuals from population 1 for trait
 125 1, \mathbf{a}_2 is the analogous vector for individuals from population 2 for trait 2, σ_1^2 and σ_2^2 are
 126 genetic variances for the two traits, σ_{12} is the genetic covariance between the traits, \mathbf{G}_{11} is a
 127 matrix with genomic relationships within population 1, \mathbf{G}_{22} is a matrix with genomic
 128 relationships within population 2, and \mathbf{G}_{12} and $\mathbf{G}_{21}(= \mathbf{G}_{12}')$ are matrices with genomic
 129 relationships between population 1 and 2.

130 To derive the definition of the genomic relationships in Equation 2, we derive the
 131 variances and covariance of the additive genetic values for the two traits. Naturally, this will
 132 result in an equation to calculate the genomic relationship matrix (\mathbf{G}) across populations to
 133 estimate (co)variances in the current populations.

134 When both populations are in Hardy-Weinberg equilibrium, allele substitution effects are
 135 independent from allele frequency, and effects of causal loci are independent from each other,
 136 the genetic variance for trait 1 can be written as $\sigma_1^2 = \sum 2p_{1j}(1-p_{1j})\sigma_{\alpha_1}^2$, where p_{1j} is the
 137 allele frequency at locus j in population 1 (Falconer and Mackay 1996). Hence, the variance
 138 of \mathbf{a}_1 is:

$$139 \quad \text{Var}(\mathbf{a}_1) = \text{Var}(\mathbf{Z}_1\boldsymbol{\alpha}_1) = \mathbf{Z}_1\mathbf{Z}_1'\sigma_{\alpha_1}^2 = \frac{\mathbf{Z}_1\mathbf{Z}_1'}{\sum 2p_{1j}(1-p_{1j})}\sigma_1^2, \quad (3)$$

140 where \mathbf{Z}_1 is a $n_I \times n_c$ matrix of centered genotypes for all individuals from population 1 (n_I)
 141 for all causal loci, and \mathbf{a}_1 is a vector of length n_c with allele substitution effects at causal loci
 142 for trait 1.

143 Similarly,

$$144 \quad \text{Var}(\mathbf{a}_2) = \frac{\mathbf{Z}_2 \mathbf{Z}_2'}{\sum 2p_{2j}(1-p_{2j})} \sigma_2^2. \quad (4)$$

145 The genetic covariance between the two traits is:

$$146 \quad \sigma_{12} = r_g \sqrt{\sigma_1^2} \sqrt{\sigma_2^2} = r_g \sqrt{\sum 2p_{1j}(1-p_{1j}) \sigma_{\alpha_1}^2} \sqrt{\sum 2p_{2j}(1-p_{2j}) \sigma_{\alpha_2}^2} =$$

$$147 \quad = \sigma_{\alpha_{12}} \sqrt{\sum 2p_{1j}(1-p_{1j})} \sqrt{\sum 2p_{2j}(1-p_{2j})}. \quad (5)$$

148 Therefore, the covariance between genetic values of population 1 and 2 is:

$$149 \quad \text{Cov}(\mathbf{a}_1, \mathbf{a}_2) = \text{Cov}(\mathbf{Z}_1 \mathbf{a}_1, \mathbf{Z}_2 \mathbf{a}_2) = \mathbf{Z}_1 \mathbf{Z}_2' \sigma_{\alpha_{12}} = \frac{\mathbf{Z}_1 \mathbf{Z}_2'}{\sqrt{\sum 2p_{1j}(1-p_{1j})} \sqrt{\sum 2p_{2j}(1-p_{2j})}} \sigma_{12}. \quad (6)$$

150 From Equation 3, 4 and 6, it follows that the genomic relationship matrix (\mathbf{G}) is:

$$151 \quad \mathbf{G} = \begin{bmatrix} \mathbf{G}_{11} & \mathbf{G}_{12} \\ \mathbf{G}_{21} & \mathbf{G}_{22} \end{bmatrix} = \begin{bmatrix} \frac{\mathbf{Z}_1 \mathbf{Z}_1'}{\sum 2p_{1j}(1-p_{1j})} & \frac{\mathbf{Z}_1 \mathbf{Z}_2'}{\sqrt{\sum 2p_{1j}(1-p_{1j})} \sqrt{\sum 2p_{2j}(1-p_{2j})}} \\ \frac{\mathbf{Z}_2 \mathbf{Z}_1'}{\sqrt{\sum 2p_{1j}(1-p_{1j})} \sqrt{\sum 2p_{2j}(1-p_{2j})}} & \frac{\mathbf{Z}_2 \mathbf{Z}_2'}{\sum 2p_{2j}(1-p_{2j})} \end{bmatrix} \quad (7)$$

153 When allele frequencies from the current population are used, \mathbf{G} from Equation 7 estimates
 154 current genetic (co)variances. Lourenco *et al.* (2016) presented a comparable \mathbf{G} matrix for
 155 combining purebred and crossbred animals. Note that the covariance of the genotypes
 156 between the populations, $\mathbf{Z}_2 \mathbf{Z}_1'$, is divided by the standard deviations of the genotypes in each
 157 population, $\sqrt{\sum 2p_{1j}(1-p_{1j})}$ and $\sqrt{\sum 2p_{2j}(1-p_{2j})}$. Therefore, the relationships in this \mathbf{G}
 158 matrix are defined as correlations between the individuals.

159 By interpreting $\sum 2p_{1j}(1-p_{1j})$, $\sum 2p_{2j}(1-p_{2j})$ and $\sqrt{\sum 2p_{1j}(1-p_{1j})}\sqrt{\sum 2p_{2j}(1-p_{2j})}$
 160 as scaling factors (i.e. k_1 , k_2 , and k_{12}) of \mathbf{G} , the variance-covariance matrix in Equation 2
 161 becomes:

$$162 \begin{bmatrix} \mathbf{G}_{11}\sigma_1^2 & \mathbf{G}_{12}\sigma_{12} \\ \mathbf{G}_{21}\sigma_{12} & \mathbf{G}_{22}\sigma_2^2 \end{bmatrix} = \begin{bmatrix} \mathbf{Z}_1\mathbf{Z}'_1 \frac{\sigma_1^2}{k_1} & \mathbf{Z}_1\mathbf{Z}'_2 \frac{\sigma_{12}}{k_{12}} \\ \mathbf{Z}_2\mathbf{Z}'_1 \frac{\sigma_{12}}{k_{12}} & \mathbf{Z}_2\mathbf{Z}'_2 \frac{\sigma_2^2}{k_2} \end{bmatrix}. \quad (8)$$

163 Equation 8 shows that the scaling factors of \mathbf{G} and the variance components are completely
 164 confounded. Therefore, other scaling factors of \mathbf{G} can be used to estimate the genetic
 165 correlation as:

$$166 \hat{r}_g = \frac{\frac{\hat{\sigma}_{12}}{k_{12}}}{\sqrt{\frac{\hat{\sigma}_1^2}{k_1}} \sqrt{\frac{\hat{\sigma}_2^2}{k_2}}} = \frac{\sqrt{k_1}\sqrt{k_2}}{k_{12}} \frac{\hat{\sigma}_{12}}{\sqrt{\hat{\sigma}_1^2}\sqrt{\hat{\sigma}_2^2}}. \quad (9)$$

167 Equation 9 shows that the genetic correlation is directly estimated from the variance
 168 components when the scaling factors of \mathbf{G} fulfil the property $k_{12} = \sqrt{k_1}\sqrt{k_2}$. When

169 $k_{12} \neq \sqrt{k_1}\sqrt{k_2}$, the correlation based on variance components should be multiplied by

170 $\left(\frac{\sqrt{k_1}\sqrt{k_2}}{k_{12}}\right)$ to correct the estimated genetic correlation. By changing the scaling factors, the

171 genetic variances change as well. When genetic variances of the current population are of
 172 interest, the within-population blocks in \mathbf{G} should be scaled as in Equation 7 (Legarra 2016).

173 Equation 8 and 9 show that the genetic correlation is estimated when the scaling factors in

174 \mathbf{G} are the same for all blocks. When all scaling factors are equal to 1, so effectively no scaling

175 factor is used, the (co)variances represent the (co)variances of the causal effects i.e., $\sigma_{\alpha_1}^2$, $\sigma_{\alpha_2}^2$

176 , and $\sigma_{\alpha_{12}}$. A disadvantage of this scaling is that elements of \mathbf{G} can become very large, which

177 can result in very small variance components that may be flagged as too small in statistical

178 software. This might be prevented by either scaling up the phenotypic variance by multiplying
179 all phenotypes by a constant, or by scaling down the elements in **G** by dividing all elements
180 by the same constant. Both scaling approaches have no influence on the genetic correlation,
181 but do affect the genetic (co)variances.

182

183 **Simulations**

184 Simulations were used to validate the results above. Two populations of 2500 individuals
185 each with phenotypes for a trait influenced by the same 15 000 loci were simulated. Allele
186 frequencies of the loci were sampled from a U-shape distribution, independently in both
187 populations. Genotypes were allocated to individuals according to Hardy-Weinberg
188 equilibrium, assuming that loci were segregating independently. Therefore, genetic
189 correlations between populations were only affected by pleiotropy and not by linkage
190 disequilibrium.

191 Allele substitution effects were sampled from a bi-variate normal distribution with means
192 zero and variances 1, and a correlation of 0.5 between allele substitution effects in both
193 populations. The allele substitution effects were multiplied with the corresponding genotypes
194 to calculate additive genetic values for individuals, assuming additive gene action.

195 Environmental effects were sampled from a normal distribution with variance ($\frac{1}{h^2} - 1$) times
196 the genetic variance, where the genetic variance was calculated across all individuals in both
197 populations. The heritability was set to 0.9, to ensure that there was sufficient power in the
198 data to estimate the (co)variances. Phenotypes were the sum of additive genetic and
199 environmental effects, and were standardized to an average of 0 and a standard deviation of
200 100. Simulations were replicated 100 times.

201 Phenotypes were analyzed in a two-trait model, using four different **G** matrices; two **G**
202 matrices derived above, and two commonly used **G** matrices for multiple populations (Chen

203 *et al.* 2013; Makgahlela *et al.* 2013). In all four methods, genotypes at causal loci were used to
 204 calculate \mathbf{G} . The methods differed in scaling factors as well as in centering of genotypes,
 205 being performed either within or across populations.

206 In the first three methods, the genotypes in \mathbf{Z} were centered within population as
 207 $g_{ijm} - 2p_{jm}$, where g_{ijm} is the allele count of individual i from population m at locus j and p_{jm}
 208 is the allele frequency at locus j in population m . The first method, \mathbf{G}_{New} , scaled \mathbf{G}
 209 following Equation 9:

$$210 \quad \mathbf{G}_{\text{New}} = \begin{bmatrix} \frac{\mathbf{Z}_1 \mathbf{Z}'_1}{\sum 2p_{1j}(1-p_{1j})} & \frac{\mathbf{Z}_1 \mathbf{Z}'_2}{\sqrt{\sum 2p_{1j}(1-p_{1j})} \sqrt{\sum 2p_{2j}(1-p_{2j})}} \\ \frac{\mathbf{Z}_2 \mathbf{Z}'_1}{\sqrt{\sum 2p_{1j}(1-p_{1j})} \sqrt{\sum 2p_{2j}(1-p_{2j})}} & \frac{\mathbf{Z}_2 \mathbf{Z}'_2}{\sum 2p_{2j}(1-p_{2j})} \end{bmatrix}.$$

211 In the second method, \mathbf{G}_{-1} , scaling factors were equal to 1:

$$212 \quad \mathbf{G}_{-1} = \begin{bmatrix} \mathbf{Z}_1 \mathbf{Z}'_1 & \mathbf{Z}_1 \mathbf{Z}'_2 \\ \mathbf{Z}_2 \mathbf{Z}'_1 & \mathbf{Z}_2 \mathbf{Z}'_2 \end{bmatrix}.$$

213 The third method, \mathbf{G}_{Chen} , calculated \mathbf{G} according to Chen *et al.* (2013):

$$214 \quad \mathbf{G}_{\text{Chen}} = \begin{bmatrix} \frac{\mathbf{Z}_1 \mathbf{Z}'_1}{\sum 2p_{1j}(1-p_{1j})} & \frac{\mathbf{Z}_1 \mathbf{Z}'_2}{\sum 2\sqrt{p_{1j}(1-p_{1j})p_{2j}(1-p_{2j})}} \\ \frac{\mathbf{Z}_2 \mathbf{Z}'_1}{\sum 2\sqrt{p_{1j}(1-p_{1j})p_{2j}(1-p_{2j})}} & \frac{\mathbf{Z}_2 \mathbf{Z}'_2}{\sum 2p_{2j}(1-p_{2j})} \end{bmatrix}.$$

215 The fourth method, $\mathbf{G}_{\text{Across}}$, used the average allele frequency across both populations
 216 instead of population-specific allele frequencies to center the genotypes (e.g., Makgahlela *et*
 217 *al.* 2013). Thus, the matrix of genotypes, denoted \mathbf{Z}^* , had elements $g_{ijm} - 2\bar{p}_j$, where \bar{p}_j is
 218 the average allele frequency across both populations at locus j . The scaling factor was the
 219 same for all blocks:

220
$$\mathbf{G_Across} = \begin{bmatrix} \frac{\mathbf{Z}_1^* \mathbf{Z}_1^{*'}}{\sum 2\bar{p}_j(1-\bar{p}_j)} & \frac{\mathbf{Z}_1^* \mathbf{Z}_2^{*'}}{\sum 2\bar{p}_j(1-\bar{p}_j)} \\ \frac{\mathbf{Z}_2^* \mathbf{Z}_1^{*'}}{\sum 2\bar{p}_j(1-\bar{p}_j)} & \frac{\mathbf{Z}_2^* \mathbf{Z}_2^{*'}}{\sum 2\bar{p}_j(1-\bar{p}_j)} \end{bmatrix}.$$

221 $\mathbf{G_New}$, $\mathbf{G_1}$ and $\mathbf{G_Across}$ fulfilled the property $k_{12} = \sqrt{k_1} \sqrt{k_2}$ to directly estimate the
 222 genetic correlation. In $\mathbf{G_Chen}$, $k_{12} \neq \sqrt{k_1} \sqrt{k_2}$ when allele frequencies in the two populations
 223 were different. Therefore, the genetic correlation estimated with $\mathbf{G_Chen}$ was multiplied by

224
$$\frac{\sqrt{k_1} \sqrt{k_2}}{k_{12}} = \frac{\sqrt{\sum 2p_{1j}(1-p_{1j})} \sqrt{\sum 2p_{2j}(1-p_{2j})}}{\sum 2\sqrt{p_{1j}(1-p_{1j})p_{2j}(1-p_{2j})}}$$
 to correct the estimate. Moreover, the current

225 populations were the base population for the within-population blocks of $\mathbf{G_New}$ and
 226 $\mathbf{G_Chen}$, so those \mathbf{G} matrices estimated the genetic variances within the current populations
 227 (Speed and Balding 2015; Legarra 2016). As explained before, the variances of $\mathbf{G_1}$
 228 represented the variances of the causal effects. For $\mathbf{G_Across}$, the base population was not
 229 clearly defined, so the interpretation of the estimated genetic variances is unclear. See
 230 supporting information for the R-script and seeds used to simulate genotypes and phenotypes
 231 and to calculate the different \mathbf{G} matrices.

232

RESULTS

233 Variance components

234 In Figure 1, the estimated genetic variance using \mathbf{G}_{New} is plotted against the simulated
235 genetic variance. This figure shows that the estimates varied only slightly around the
236 simulated values. This shows that \mathbf{G}_{New} unbiasedly estimated the genetic variance in the
237 current populations.

238 As expected, \mathbf{G}_{New} and \mathbf{G}_{Chen} estimated the same genetic variances (Figure 2 and 3).
239 The variances of \mathbf{G}_{1} represented the variances of the causal effects. By multiplying those
240 variances by $\sum 2p_{jm}(1-p_{jm})$ for population m , genetic variances identical to \mathbf{G}_{New} and
241 \mathbf{G}_{Chen} were obtained. The genetic variance estimated with $\mathbf{G}_{\text{Across}}$ was approximately a
242 factor 1.5 higher than the genetic variance estimated with \mathbf{G}_{New} and \mathbf{G}_{Chen} . Also the
243 scaling factors k_1 and k_2 were approximately a factor 1.5 higher. Hence, when multiplying the
244 variances estimated with $\mathbf{G}_{\text{Across}}$ by the ratio in scaling factors, estimates became identical
245 to those with \mathbf{G}_{New} and \mathbf{G}_{Chen} . So, the difference in estimated variances between
246 methods was completely explained by the difference in scaling factors, while centering
247 genotypes within or across populations had no effect on estimated variances. Estimated
248 residual variances were exactly the same for the four different \mathbf{G} matrices.

249

250 Genetic correlation

251 Despite differences in (co)variance estimates, \mathbf{G}_{New} , \mathbf{G}_{1} , and $\mathbf{G}_{\text{Across}}$ yielded the
252 same estimated genetic correlation (Figure 4) which was an unbiased estimate of the
253 simulated genetic correlation (Figure 5). This is because differences in genetic covariances
254 among models were compensated by corresponding differences in genetic variances. The

255 genetic correlation estimated using \mathbf{G}_{Chen} was ~20% lower. When multiplying this estimate

256 by $\frac{\sqrt{k_1}\sqrt{k_2}}{k_{12}}=1.23$, the genetic correlation became identical to the other three methods.

257

258

DISCUSSION

259 The aim of this paper was to define a multi-population genomic relationship matrix to
260 estimate current genetic variances within and genetic correlations between populations. We
261 derived a genomic relationship matrix, \mathbf{G}_{New} , that yields unbiased estimates of current
262 genetic variances, covariances and correlations. Moreover, we showed the required property
263 for other genomic relationship matrices to estimate the genetic correlation between
264 populations, even though estimated variance components are not necessarily related to the
265 current populations.

266

267 **Methods to calculate the genomic relationship matrix**

268 From the four methods used in this paper to calculate \mathbf{G} , \mathbf{G}_{New} was the only matrix
269 correctly estimating both current genetic variances as well as genetic correlations. \mathbf{G}_{Chen}
270 also estimated current genetic variances, but the estimated genetic correlation had to be
271 multiplied by $\frac{\sqrt{k_1}\sqrt{k_2}}{k_{12}}$. \mathbf{G}_1 estimated the correct genetic correlation, but estimated the
272 variance of causal effects instead of the genetic variance. Although the base population in
273 $\mathbf{G}_{\text{Across}}$ was not well defined, genetic correlations were correctly estimated but there was
274 no clear interpretation of the estimated genetic variances. Results also showed that genetic
275 variances were not affected by centering the allele count, as shown before by Strandén and
276 Christensen (2011).

277 Table 1 gives an overview of the most frequently used methods to calculate \mathbf{G} across
278 multiple populations, with scaling factors and correction factors for the estimated genetic
279 correlation. \mathbf{G}_{New} , \mathbf{G}_1 , $\mathbf{G}_{\text{Across}}$, and the method described by Erbe *et al.* (2012) directly
280 estimate the correct genetic correlation. The \mathbf{G}_{Chen} method does not directly estimate the
281 genetic correlation, but the estimate can be corrected using the scaling factors. Those five

282 methods all assume that allele substitution effects are independent of allele frequency, similar
283 to method 1 of VanRaden (2008). This is in contrast to another regularly used method, namely
284 method 2 of VanRaden (2008), also described by Yang (2010). This method yields a valid
285 relationship matrix only when the average effect at a locus is proportional to the reciprocal of
286 the square root of expected heterozygosity at that locus (Appendix, Equation A8). So, this
287 method assumes that marker effects are determined by their allele frequency, with larger
288 effects for rarer alleles. For a trait determined by relatively few genes and undergoing
289 directional selection, this assumption may be plausible, since selection acts stronger on causal
290 loci with a larger effect (Haldane 1924; Wright 1931, 1937). It is, however, a very strong
291 assumption in general. Many traits may experience only weak selection, and/or are
292 determined by many genes. In those cases, allele frequency distribution is determined mainly
293 by the interplay of mutation and drift, and a direct relationship between effect size and allele
294 frequency is not expected. Therefore, the assumption of independence between allele
295 frequency and allele substitution effects seems more realistic for most traits. Moreover, when
296 allele substitution effects would depend on allele frequency, effects for exactly the same trait
297 would differ between populations when allele frequencies differ. This makes the
298 interpretation of a genetic correlation estimated using method 2 of VanRaden (2008) rather
299 difficult. Therefore, we advise to use \mathbf{G} matrices based on method 1 instead of method 2 of
300 VanRaden (2008), especially when multiple populations are considered.

301 In this paper, we assumed that causal loci were known and were used to calculate \mathbf{G} . In
302 this way, differences in linkage disequilibrium (LD) between markers and causal loci across
303 populations did not affect the results and all genetic variance was explained by \mathbf{G} . When
304 genomic markers are used to calculate \mathbf{G} , differences in LD can affect the results, since the
305 LD pattern is known to differ across populations in humans (Sawyer *et al.* 2005) as well as in
306 livestock (e.g., Heifetz *et al.* 2005; Gautier *et al.* 2007; Veroneze *et al.* 2013). This difference

307 in LD is likely to affect the estimated genetic correlation, since it reduces the correlation of
308 marker effects (Gianola *et al.* 2015). Moreover, markers might not explain all genetic
309 variance when there is no complete LD between a causal locus and at least one marker (e.g.,
310 Yang *et al.* 2010; Daetwyler *et al.* 2013). This can affect the estimated genetic correlation
311 when the variance explained by the markers shows either a higher or lower genetic correlation
312 than the part not explained (Bulik-Sullivan *et al.* 2015). Therefore, it is difficult to predict the
313 effect of not explaining all genetic variance by markers on the estimated genetic correlation.
314 In a follow-up study, we will investigate the effect of using marker genotypes on the
315 estimated genetic correlation between populations.

316

317 **Other approaches to estimate the genetic correlation between populations**

318 We focused on using genomic relationships in a multi-trait model to estimate genetic
319 correlations between populations. Genetic correlations can also be estimated using summary
320 statistics of genome-wide association studies (GWAS; Bulik-Sullivan *et al.* 2015; Brown *et*
321 *al.* 2016) or using random regression on genotypes (Sørensen *et al.* 2012; Krag *et al.* 2013).
322 The method based on summary statistics of GWAS combines information from different
323 studies and weights estimated marker effects by LD overlap and corresponding z score (Bulik-
324 Sullivan *et al.* 2015; Brown *et al.* 2016). This method is beneficial when the costs of
325 collecting enough data are high and data sharing is not possible. It is, however, not known
326 whether this method estimates the correct genetic correlation. The method using random
327 regression on genotypes is equivalent to the multi-trait GREML method used in this study,
328 since both estimate the same additive genetic values when the genotypes are centered and
329 scaled in the same way (Habier *et al.* 2007; VanRaden 2008; Goddard 2009). Variance
330 components estimated with random regression on marker genotypes represent variances of
331 marker effects (Meuwissen *et al.* 2001), similar to \mathbf{G}_1 , when the same centered genotypes

332 are used as input. Hence, random regression on centered genotypes can also be used to
333 estimate genetic correlations between populations. When genotypes for the random regression
334 are centered and scaled, the estimated genetic correlation becomes equal to the estimated
335 genetic correlation using \mathbf{G} based on method 2 of VanRaden (VanRaden 2008; Yang *et al.*
336 2010). Therefore, the interpretation of this estimated genetic correlation remains unclear as
337 well.

338

339 **Importance of the genetic correlation between populations**

340 The genetic correlation between populations is an important parameter for genomic
341 prediction, since it determines the usefulness of combining information from multiple
342 populations. A low genetic correlation means that it is very unlikely that combining
343 populations will increase the accuracy of estimated genetic values. Therefore, the genetic
344 correlation partly determines the accuracy of across- or multi-population genomic prediction.
345 For predicting the accuracy in those scenarios, an accurate estimation of genetic correlations
346 is essential (Wientjes *et al.* 2015; Wientjes *et al.* 2016). For predicting response to selection,
347 both the accuracy as well as current genetic variances are needed (Falconer and Mackay
348 1996). Even though the accuracy of estimated genetic values is quite consistent across
349 methods for calculating \mathbf{G} (Makgahlela *et al.* 2013, 2014; Lourenco *et al.* 2016), for
350 estimating genetic (co)variances and correlations it is important to use the \mathbf{G}_{New} matrix.

351

352 **Genetic correlation versus genic correlation**

353 The genetic correlation is defined based on additive genetic (co)variances. Under selection,
354 however, additive genetic (co)variances change over generations, since selection creates
355 transient gametic phase disequilibrium (i.e., correlations between allele substitution effects at
356 different loci). This process is also known as the Bulmer effect (Bulmer 1971). Therefore,

357 genetic (co)variances and correlations depend not only on the genetic background of the traits,
358 but also on transient processes like the type and intensity of selection. Apart from additive
359 genetic (co)variances, quantitative genetics also describes genic (co)variances (e.g., Bulmer
360 1980; Bulmer 1989), defined as the additive genetic (co)variance in the absence of gametic
361 phase disequilibrium. In contrast to genetic variances, genic variances are independent from
362 selection and are always equal to twice the Mendelian sampling variance (Hill 2014). In
363 analogy to genic (co)variances, genic correlations can be defined as well. We believe that
364 genic correlations are more relevant than additive genetic correlations, since genic
365 correlations are not influenced by transient processes and, therefore, more constant across
366 generations.

367 In our simulation study, allele substitution effects were randomly sampled, so no transient
368 gametic phase disequilibrium was present and genic (co)variances were equal to the additive
369 genetic (co)variances. In all situations, genic variances can be estimated when the base
370 population of the relationship matrix is unselected and phenotypic records on which selection
371 decisions are based are available (Henderson 1985). It is also shown that even when
372 phenotypic records from the base population are absent, the genic variance can be estimated
373 when phenotypic records for several generations are available and the base population is
374 unselected (Henderson 1985; Van der Werf and de Boer 1990). It can be expected that as long
375 as several generations of phenotypic data is available in combination with the relationships
376 between all those individuals, variances are corrected for selection and effectively genic
377 variances are estimated. Therefore, genic correlations can likely be calculated using \mathbf{G}_{New} ,
378 provided that data is available for several generations.

379

380 **Conclusion**

381 The properties of the genomic relationship matrix affect estimates of genetic variances
382 within as well as genetic correlations between populations. For estimating current genetic
383 variances, allele frequencies of the current population should be used to calculate
384 relationships within that population. For estimating genetic correlations between populations,
385 scaling factors of the different blocks of the relationship matrix, based on method 1 of
386 VanRaden (2008), should fulfill the property $k_{12} = \sqrt{k_1} \sqrt{k_2}$. When this property is not
387 fulfilled, the estimated genetic correlation can be corrected by multiplying the estimate by
388 $\frac{\sqrt{k_1} \sqrt{k_2}}{k_{12}}$. In this study we present a genomic relationship matrix, **G**_New, which directly
389 results in current genetic variances as well as genetic correlations between populations.

390

391

392

393

ACKNOWLEDGMENTS

394 This research is supported by the Netherlands Organisation of Scientific Research (NWO)

395 and the Breed4Food consortium partners Cobb Europe, CRV, Hendrix Genetics and

396 Topigs Norsvin.

397

398

399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423

LITERATURE CITED

- Bohren, B. B., W. G. Hill and A. Robertson, 1966 Some observations on asymmetrical correlated responses to selection. *Genet. Res.* 7: 44-57.
- Brown, B. C., C. J. Ye, A. L. Price and N. Zaitlen, 2016 Transethnic genetic-correlation estimates from summary statistics. *Am. J. Hum. Genet.* 99: 76-88.
- Bulik-Sullivan, B., H. K. Finucane, V. Anttila, A. Gusev, F. R. Day, *et al.*, 2015 An atlas of genetic correlations across human diseases and traits. *Nat. Genet.* 47: 1236-1241.
- Bulmer, M., 1989 Maintenance of genetic variability by mutation-selection balance: a child's guide through the jungle. *Genome* 31: 761-767.
- Bulmer, M. G., 1971 The effect of selection on genetic variability. *Am. Nat.* 105: 201-211.
- Bulmer, M. G., 1980 *The mathematical theory of quantitative genetics*. Clarendon Press., Oxford.
- Chen, L., F. Schenkel, M. Vinsky, D. Crews and C. Li, 2013 Accuracy of predicting genomic breeding values for residual feed intake in Angus and Charolais beef cattle. *J. Anim. Sci.* 91: 4669-4678.
- Daetwyler, H. D., M. P. L. Calus, R. Pong-Wong, G. De los Campos and J. M. Hickey, 2013 Genomic prediction in animals and plants: Simulation of data, validation, reporting, and benchmarking. *Genetics* 193: 347-365.
- De Candia, T. R., S. H. Lee, J. Yang, B. L. Browning, P. V. Gejman, *et al.*, 2013 Additive genetic variation in schizophrenia risk is shared by populations of African and European descent. *Am. J. Hum. Genet.* 93: 463-470.
- Erbe, M., B. J. Hayes, L. K. Matukumalli, S. Goswami, P. J. Bowman, *et al.*, 2012 Improving accuracy of genomic predictions within and between dairy cattle breeds with imputed high-density single nucleotide polymorphism panels. *J. Dairy Sci.* 95: 4114-4129.

- 424 Falconer, D. S., 1952 The problem of environment and selection. *Amer. Nat.* 86: 293-298.
- 425 Falconer, D. S. and T. F. C. Mackay, 1996 *Introduction to quantitative genetics*. Pearson
426 Education Limited, Harlow.
- 427 Fisher, R. A., 1918 The Correlation between relatives on the supposition of Mendelian
428 inheritance. *Trans. Roy. Soc. Edinburgh* 52: 399-433.
- 429 Fisher, R. A., 1930 *The genetical theory of natural selection: a complete variorum edition*.
430 Oxford University Press.
- 431 Gautier, M., T. Faraut, K. Moazami-Goudarzi, V. Navratil, M. Foglio, *et al.*, 2007 Genetic
432 and haplotypic structure in 14 European and African cattle breeds. *Genetics* 177: 1059-
433 1070.
- 434 Gianola, D., G. De los Campos, M. A. Toro, H. Naya, C.-C. Schön, *et al.*, 2015 Do molecular
435 markers inform about pleiotropy? *Genetics* 201: 23-29.
- 436 Goddard, M. E., 2009 Genomic selection: Prediction of accuracy and maximisation of long
437 term response. *Genetica* 136: 245-257.
- 438 Habier, D., R. L. Fernando and J. C. M. Dekkers, 2007 The impact of genetic relationship
439 information on genome-assisted breeding values. *Genetics* 177: 2389-2397.
- 440 Haldane, J. B. S., 1924 A mathematical theory of natural and artificial selection—I. *Trans.*
441 *Camb. Phil. Soc.* 23: 19-41.
- 442 Harris, B. L. and D. L. Johnson, 2010 Genomic predictions for New Zealand dairy bulls and
443 integration with national genetic evaluation. *J. Dairy Sci.* 93: 1243-1252.
- 444 Hayes, B. J., P. M. Visscher and M. E. Goddard, 2009 Increased accuracy of artificial
445 selection by using the realized relationship matrix. *Genet. Res.* 91: 47-60.
- 446 Heifetz, E. M., J. E. Fulton, N. O'Sullivan, H. Zhao, J. C. M. Dekkers, *et al.*, 2005 Extent and
447 consistency across generations of linkage disequilibrium in commercial layer chicken
448 breeding populations. *Genetics* 171: 1173-1181.

- 449 Henderson, C. R., 1985 Best linear unbiased prediction using relationship matrices derived
450 from selected base populations. *J. Dairy Sci.* 68: 443-448.
- 451 Hill, W. G., 2014 Applications of population genetics to animal breeding, from Wright, Fisher
452 and Lush to genomic prediction. *Genetics* 196: 1-16.
- 453 Karoui, S., M. Carabaño, C. Díaz and A. Legarra, 2012 Joint genomic evaluation of French
454 dairy cattle breeds using multiple-trait models. *Genet. Sel. Evol.* 44: 39.
- 455 Krag, K., N. A. Poulsen, M. K. Larsen, L. B. Larsen, L. L. Janss, *et al.*, 2013 Genetic
456 parameters for milk fatty acids in Danish Holstein cattle based on SNP markers using a
457 Bayesian approach. *BMC Genet.* 14: 79.
- 458 Legarra, A., 2016 Comparing estimates of genetic variance across different relationship
459 models. *Theor. Popul. Biol.* 107: 26-30.
- 460 Lehermeier, C., C.-C. Schön and G. De los Campos, 2015 Assessment of genetic
461 heterogeneity in structured plant populations using multivariate whole-genome regression
462 models. *Genetics* 201: 323-337.
- 463 Lourenco, D. A. L., S. Tsuruta, B. O. Fragomeni, C. Y. Chen, W. O. Herring, *et al.*, 2016
464 Crossbreed evaluations in single-step genomic best linear unbiased predictor using
465 adjusted realized relationship matrices. *J. Anim. Sci.* 94: 909-919.
- 466 Makgahlela, M. L., I. Strandén, U. S. Nielsen, M. J. Sillanpää and E. A. Mäntysaari, 2013 The
467 estimation of genomic relationships using breedwise allele frequencies among animals in
468 multibreed populations. *J. Dairy Sci.* 96: 5364-5375.
- 469 Makgahlela, M. L., I. Strandén, U. S. Nielsen, M. J. Sillanpää and E. A. Mäntysaari, 2014
470 Using the unified relationship matrix adjusted by breed-wise allele frequencies in genomic
471 evaluation of a multibreed population. *J. Dairy Sci.* 97: 1117-1127.
- 472 Meuwissen, T. H. E., B. J. Hayes and M. E. Goddard, 2001 Prediction of total genetic value
473 using genome-wide dense marker maps. *Genetics* 157: 1819-1829.

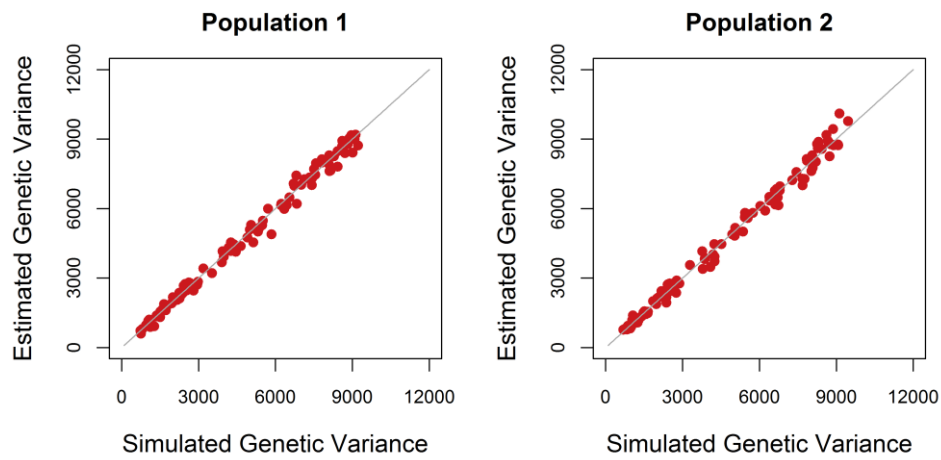
- 474 Olson, K. M., P. M. VanRaden and M. E. Tooker, 2012 Multibreed genomic evaluations
475 using purebred Holsteins, Jerseys, and Brown Swiss. *J. Dairy Sci.* 95: 5378-5383.
- 476 Powell, J. E., P. M. Visscher and M. E. Goddard, 2010 Reconciling the analysis of IBD and
477 IBS in complex trait studies. *Nat. Rev. Gen.* 11: 800-805.
- 478 Sawyer, S. L., N. Mukherjee, A. J. Pakstis, L. Feuk, J. R. Kidd, *et al.*, 2005 Linkage
479 disequilibrium patterns vary substantially among populations. *Europ. J. Hum. Genet.* 13:
480 677-686.
- 481 Sørensen, L. P., L. Janss, P. Madsen, T. Mark and M. S. Lund, 2012 Estimation of
482 (co)variances for genomic regions of flexible sizes: application to complex infectious
483 udder diseases in dairy cattle. *Genet. Sel. Evol.* 44: 18.
- 484 Speed, D. and D. J. Balding, 2015 Relatedness in the post-genomic era: is it still useful? *Nat.*
485 *Rev. Genet.* 16: 33-44.
- 486 Strandén, I. and O. F. Christensen, 2011 Allele coding in genomic evaluation. *Genet. Sel.*
487 *Evol.* 43: 25.
- 488 Van der Werf, J. and I. de Boer, 1990 Estimation of additive genetic variance when base
489 populations are selected. *J. Anim. Sci.* 68: 3124-3132.
- 490 VanRaden, P. M., 2008 Efficient methods to compute genomic predictions. *J. Dairy Sci.* 91:
491 4414-4423.
- 492 Veroneze, R., P. S. Lopes, S. E. F. Guimarães, F. F. Silva, M. S. Lopes, *et al.*, 2013 Linkage
493 disequilibrium and haplotype block structure in six commercial pig lines. *J. Anim. Sci.* 91:
494 3493-3501.
- 495 Visscher, P. M., G. Hemani, A. A. E. Vinkhuyzen, G.-B. Chen, S. H. Lee, *et al.*, 2014
496 Statistical power to detect genetic (co)variance of complex traits using SNP data in
497 unrelated samples. *PLoS Genet* 10: e1004269.

- 498 Wientjes, Y. C. J., R. F. Veerkamp, P. Bijma, H. Bovenhuis, C. Schrooten, *et al.*, 2015
499 Empirical and deterministic accuracies of across-population genomic prediction. *Genet.*
500 *Sel. Evol.* 47: 5.
- 501 Wientjes, Y. C. J., P. Bijma, R. F. Veerkamp and M. P. L. Calus, 2016 An equation to predict
502 the accuracy of genomic values by combining data from multiple traits, breeds, lines, or
503 environments. *Genetics* 202: 799-823.
- 504 Wright, S., 1931 Evolution in Mendelian populations. *Genetics* 16: 97-159.
- 505 Wright, S., 1937 The distribution of gene frequencies in populations. *Proc. Nat. Acad. Sci.*
506 *USA* 23: 307-320.
- 507 Yang, J., B. Benyamin, B. P. McEvoy, S. Gordon, A. K. Henders, *et al.*, 2010 Common SNPs
508 explain a large proportion of the heritability for human height. *Nat. Genet.* 42: 565-569.
- 509 Yang, J., A. Bakshi, Z. Zhu, G. Hemani, A. A. E. Vinkhuyzen, *et al.*, 2015 Genetic variance
510 estimation with imputed variants finds negligible missing heritability for human height and
511 body mass index. *Nat. Genet.* 47: 1114-1120.
- 512

513

FIGURES

514



515

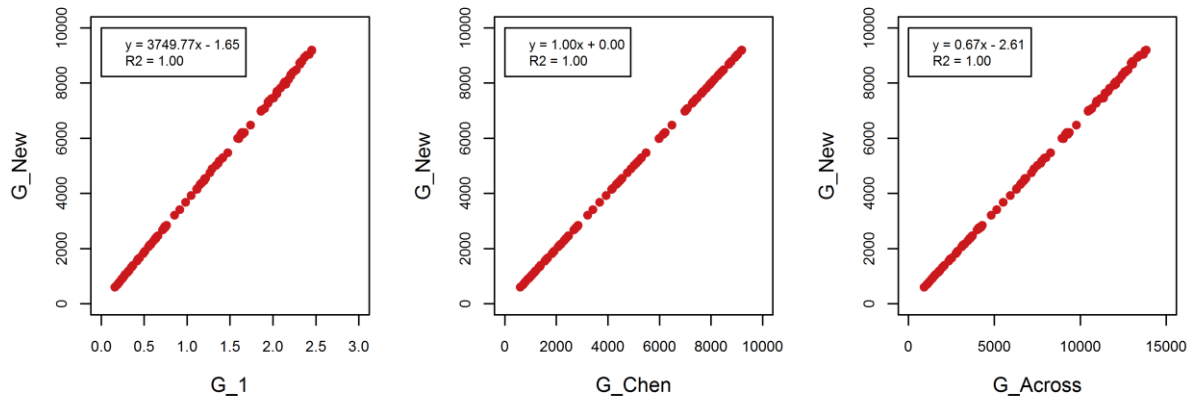
516 **Figure 1 – Estimated versus simulated genetic variance.** The estimated genetic variance in
517 both populations in each of the 100 replicates using the genomic relationship matrix derived
518 in this study (G_{New}) versus the simulated genetic variance. The grey line represents the line
519 $y=x$.

520

521

522

Genetic variance - Population 1



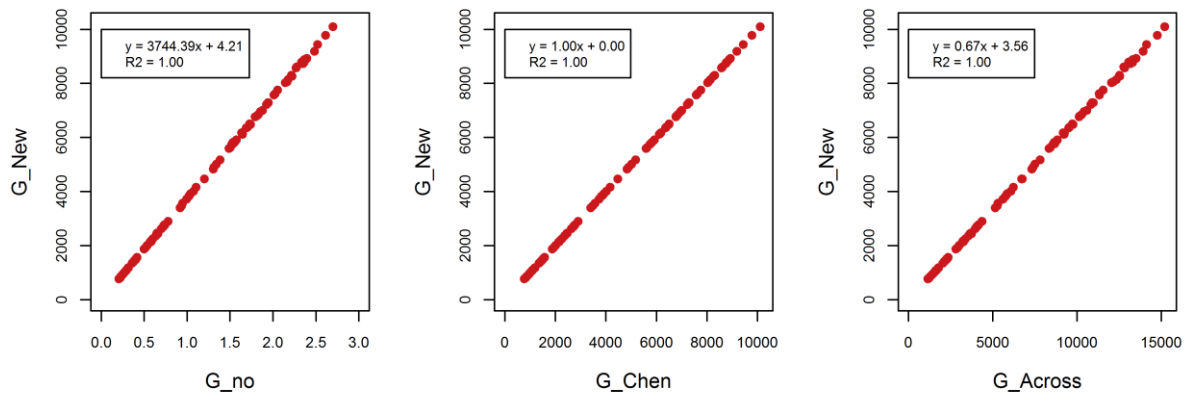
523

524 **Figure 2 – Estimated genetic variance in population 1.** The estimated genetic variance in
525 population 1 in each of the 100 replicates using the genomic relationship matrix derived in
526 this study (G_{New}) versus the estimated genetic variance using population-specific allele
527 frequencies and either a genomic relationship matrix without scaling factors (G_1) or based
528 on the method of Chen *et al.* (2013; G_{Chen}), or using allele frequencies across populations
529 (G_{Across}).

530

531

Genetic variance - Population 2



532

533 **Figure 3 – Estimated genetic variance in population 2.** The estimated genetic variance in

534 population 2 in each of the 100 replicates using the genomic relationship matrix derived in

535 this study (G_New) versus the estimated genetic variance using population-specific allele

536 frequencies and either a genomic relationship matrix without scaling factors (G_1) or based

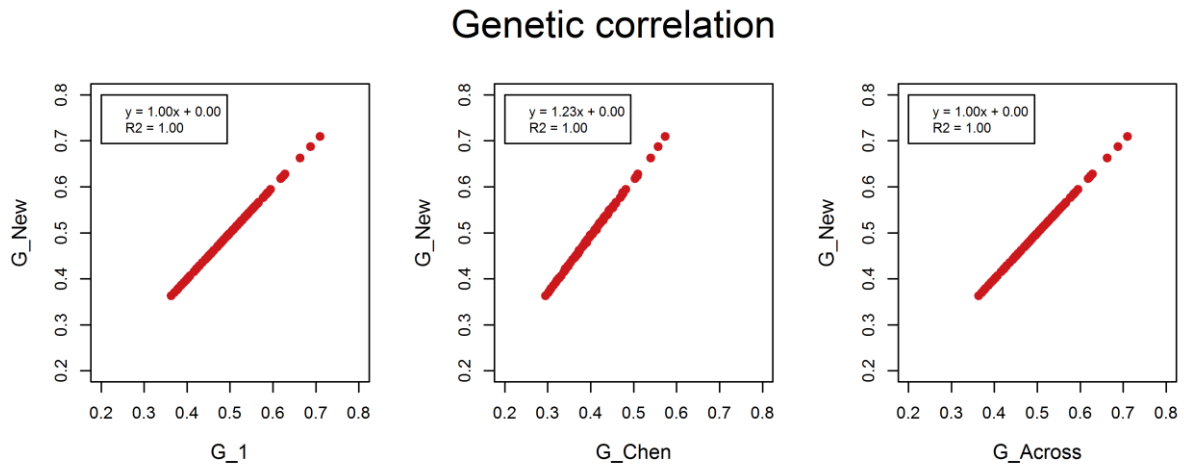
537 on the method of Chen *et al.* (2013; G_Chen), or using allele frequencies across populations

538 (G_Across).

539

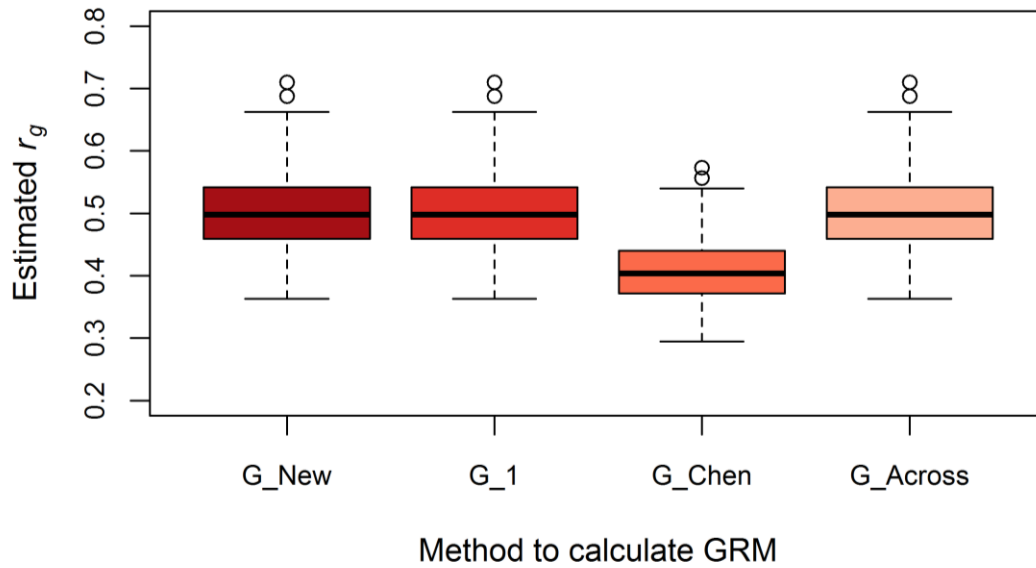
540

541



542
543 **Figure 4 – Estimated genetic correlation between population 1 and 2.** The estimated
544 genetic correlation between population 1 and 2 in each of the 100 replicates using the
545 genomic relationship derived in this study (**G_New**) versus the estimated genetic correlation
546 using population-specific allele frequencies and either a genomic relationship matrix without
547 scaling factors (**G_1**), based on the method of Chen *et al.* (2013; **G_Chen**), or using allele
548 frequencies across populations (**G_Across**).

549



550

551 **Figure 5 – Boxplot of the estimated genetic correlation using different methods to**
552 **calculate the genomic relationship matrix.** The estimated genetic correlation between
553 population 1 and 2 in each of the 100 replicates using the genomic relationship matrix derived
554 in this study (**G_New**), using population-specific allele frequencies and either a genomic
555 relationship matrix without scaling factors (**G_1**), or based on the method of Chen *et al.*
556 (2013; **G_Chen**), or using allele frequencies across populations (**G_Across**). The simulated
557 genetic correlation was 0.5.

558 **Table 1** – Overview of frequently used method to calculate **G** across populations with scaling and correction factors.

Method of calculating G ^a	Described by	Used scaling factors of the different blocks in G ^b			Correction factor to correct the genetic correlation
		k_1^c	k_2^c	k_{12}^c	
G_New	This study	$\sum 2p_{1i}(1-p_{1i})$	$\sum 2p_{2i}(1-p_{2i})$	$\sqrt{\sum 2p_{1i}(1-p_{1i})}\sqrt{\sum 2p_{2i}(1-p_{2i})}$	Not needed
G_1	This study	1	1	1	Not needed
G_Chen	Chen <i>et al.</i> (2013)	$\sum 2p_{1i}(1-p_{1i})$	$\sum 2p_{2i}(1-p_{2i})$	$\sum 2\sqrt{p_{1i}(1-p_{1i})p_{2i}(1-p_{2i})}$	$\frac{\sqrt{\sum 2p_{1i}(1-p_{1i})}\sqrt{\sum 2p_{2i}(1-p_{2i})}}{\sum 2\sqrt{p_{1i}(1-p_{1i})p_{2i}(1-p_{2i})}}$
G_Across	VanRaden (2008)/ Makgahlela <i>et al.</i> (2013)	$\sum 2\bar{p}_i(1-\bar{p}_i)$	$\sum 2\bar{p}_i(1-\bar{p}_i)$	$\sum 2\bar{p}_i(1-\bar{p}_i)$	Not needed
Erbe	Erbe <i>et al.</i> (2012)	$\sum 2p_i^*(1-p_i^*)$	$\sum 2p_i^*(1-p_i^*)$	$\sum 2p_i^*(1-p_i^*)$	Not needed
VanRaden2/ Yang	VanRaden (2008); Yang <i>et al.</i> (2010)	Nr. of markers ^d	Nr. of markers ^d	Nr. of markers ^d	Unknown

559 ^a Methods were compared assuming that no adjustment for inbreeding or regression back to the pedigree relationship matrix was performed.

560 ^b k_1 is the scaling factor of the block containing relationships in population 1, k_2 is the scaling factor of the block containing relationships in
561 population 2, and k_{12} is the scaling factor of the block containing relationship between population 1 and 2.

562 ^c p_{1i} is the allele frequency in population 1, p_{2i} is the allele frequency in population 2, \bar{p}_i is the average allele frequency across populations, p_i^*
563 is the estimated allele frequency when the populations separated.

564 ^d Per marker i , genotypes are scaled by $\sqrt{2p_i(1-p_i)}$.

565

566

567

APPENDIX

568 The \mathbf{G} matrix based on method 2 of VanRaden (2008) and Yang *et al.* (2010), \mathbf{G}_{VR2} ,
 569 weights markers by the reciprocal of the square root of the variance of its genotypes. In this
 570 Appendix, it is shown that this is only correct under the assumption that the variance of the
 571 average effect (α) at a locus, say l , is inversely proportional to expected heterozygosity at that
 572 locus,

$$573 \quad \sigma_{\alpha_l}^2 = \frac{c}{2p_l(1-p_l)}, \quad (\text{A1})$$

574 where c is a constant, and p_l the allele frequency at locus l .

575 Consider the single-trait mixed model $\mathbf{y} = \mathbf{Xb} + \mathbf{Za} + \mathbf{e}$, where \mathbf{a} is the vector of random
 576 additive genetic effects, with $\text{var}(\mathbf{a}) = \mathbf{G}\sigma_A^2$. This mixed model is valid only when $\mathbf{G}\sigma_A^2$
 577 indeed represents the covariances between additive genetic effects (A) of individuals. This
 578 requires that

$$579 \quad \mathbf{G}_{ij} = \text{cov}(A_i, A_j) / \text{var}(A), \quad (\text{A2})$$

580 where i and j are individuals.

581 By definition, the additive genetic effect of an individual is the sum of the average effects
 582 at its loci, weighted by the centred allele count (Fisher 1918; Falconer and Mackay 1996),

$$583 \quad A_i = \sum_l (x_{il} - 2p_l)\alpha_l, \quad (\text{A3})$$

584 where x_{il} is the allele count of individual i at locus l , taking values 0, 1 or 2. Thus

$$585 \quad \text{cov}(A_i, A_j) = \text{cov} \left[\sum_l (x_{il} - 2p_l)\alpha_l, \sum_l (x_{jl} - 2p_l)\alpha_l \right]. \quad (\text{A4})$$

586 For the genic covariance, the $(x_{il} - 2p_l)\alpha_l$ terms are independent between loci by
 587 definition (Bulmer 1971), so that the covariance reduces to

$$588 \quad \text{cov}(A_i, A_j) = \sum_l (x_{il} - 2p_l)(x_{jl} - 2p_l)\sigma_{\alpha_l}^2. \quad (\text{A5})$$

589 Substituting the relationship between average effects and allele frequency given by
590 Equation A1 yields

$$591 \quad \text{cov}(A_i, A_j) = c \sum_l \left[\frac{(x_{il} - 2p_l)(x_{jl} - 2p_l)}{2p_l(1-p_l)} \right]. \quad (\text{A6})$$

592 Analogously, the genic variance equals

$$593 \quad \text{var}(A) = \sum_l 2p_l(1-p_l)\sigma_{\alpha_l}^2 = \sum_l c = n_l c,$$

594 where n_l is the number of loci. Finally, from Equation A2,

$$595 \quad \mathbf{G}_{ij} = \text{cov}(A_i, A_j) / \text{var}(A) = \frac{1}{n_l} \sum_l \left[\frac{(x_{il} - 2p_l)(x_{jl} - 2p_l)}{2p_l(1-p_l)} \right], \quad (\text{A7})$$

596 which is $\mathbf{G_VR2}$. Thus obtaining $\mathbf{G_VR2}$ requires Equation A1.

597 Hence, $\mathbf{G_VR2}$ is valid under the assumption that the magnitude of the average effect at a
598 locus is proportional to the reciprocal of the square root of expected heterozygosity at that
599 locus,

$$600 \quad \alpha_l \propto \frac{1}{\sqrt{2p_l(1-p_l)}}. \quad (\text{A8})$$

601 Equation A7 shows that elements of $\mathbf{G_VR2}$ are the genome-wide average of the
602 correlations at individual loci; the term in square-brackets is the correlation between additive
603 genetic effects at locus l , and the sum of these terms is divided by the number of loci. Thus
604 $\mathbf{G_VR2}$ may have been motivated as the genome-wide average of relationships at individual
605 loci.

606 However, relatedness refers to the correlation between the total additive genetic effects of
607 individuals (Equation A2), which are sums of additive genetic effects at individual loci. In
608 general, the correlation between sums does not equal the average correlation between
609 components of the sums,

610
$$\mathbf{G}_{ij} \neq \frac{1}{n_l} \sum_l \mathbf{G}_{ijl} \quad (\text{A9})$$

611 but is defined as the ratio of the covariance and variance of the sum,

612
$$\mathbf{G}_{ij} = \text{cov}(A_i, A_j) / \sigma_A^2. \quad (\text{A10})$$

613 Equations A9 and A10 are only equal to each other under the assumption given in Equation

614 A1.

615

616

617