

Supplementary Methods

Sections 1 and 2 give a detailed account of the model and of the analytical derivations presented in the main manuscript. All the results presented in Section 1 and 2.1 are well established in the literature; we have tried to provide detailed derivations based on properties of the multivariate normal density. Readers interested on further details are referred to [1,2]. To the best of our knowledge, the results presented on section 2.2., prediction R-squared under imperfect LD between markers and genotypes at causal loci, are novel. Finally **Section 3** gives a complete account of the Bayesian implementation used to fit models.

1. Genomic BLUP (Best Linear Unbiased Predictor)

Consider a linear regression on marker covariates of the form,

$$y_i = \sum_{l=1}^p w_{il} \beta_l + \varepsilon_i, \quad (\text{S.1})$$

where, y_i represents a phenotypic measure taken on the i^{th} individual in the sample ($i=1, \dots, n$), w_{il} are marker covariates ($l=1, \dots, p$ markers), β_l are marker effects and ε_i are model residuals. Here, we assume that phenotypes are centered, $\sum_{i=1}^n y_i = 0$, and that marker covariates are centered and standardized to unit variance, $w_{il} = \frac{(x_{il} - 2\theta_l)}{\sqrt{2\theta_l(1-\theta_l)}}$, where $x_{il} \in (0,1,2)$ counts the number of copies of one of the alleles observed at the l^{th} locus of the i^{th} individual and θ_l is an estimate of the frequency of the allele coded as one at the l^{th} locus. Centering and standardization are not strictly needed for the arguments we outline to hold, but the presentation is greatly facilitated. Stacking all the equations for individuals 1 to n , we have:

$$\mathbf{y} = \mathbf{W}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \quad (\text{S.2})$$

where $\mathbf{y} = \{y_i\}$, $\boldsymbol{\beta} = \{\beta_l\}$ and $\boldsymbol{\varepsilon} = \{\varepsilon_i\}$ represent vectors of phenotypes, marker effects and model residuals and $\mathbf{W} = \{w_{il}\}$ is a matrix of marker covariates. In a standard Bayesian Gaussian Regression marker effects are assumed to be IID (independent and identically distributed) draws from a normal

25 density $\beta_l \stackrel{iid}{\sim} N(0, \sigma_\beta^2)$, and model residuals are also assumed to be IID normal, $\varepsilon_i \stackrel{iid}{\sim} N(0, \sigma_\varepsilon^2)$,
 26 independent of marker effects. Here, σ_β^2 and σ_ε^2 represent the prior variance of marker effects and the
 27 variance of model residuals.

28 **Model.** The linear score $u_i = \sum_{l=1}^p w_{il}\beta_l$, the ‘genomic value’, represents the expected value of
 29 the i^{th} phenotype given marker genotypes and marker effects. Replacing this in the data equation (S.1) we
 30 arrive at the following random effects model $y_i = u_i + \varepsilon_i$, or in matrix notation

$$31 \quad \mathbf{y} = \mathbf{u} + \boldsymbol{\varepsilon} \quad (S.3)$$

32 where, $\mathbf{u} = (u_1, \dots, u_n)'$. Following standard properties of the multivariate normal density it can be shown
 33 that the joint density of \mathbf{u} is multivariate normal, with mean equal to zero and variance-covariance matrix
 34 proportional to $\mathbf{G} = p^{-1}\mathbf{W}\mathbf{W}'$; therefore, $\mathbf{u} \sim N(\mathbf{0}, \mathbf{G}\sigma_u^2)$, with $\sigma_u^2 = p\sigma_\beta^2$.

35 Collecting assumptions, the joint distribution of phenotypes, genetic values and model residuals is
 36 given by:

$$37 \quad \begin{bmatrix} \mathbf{u} \\ \boldsymbol{\varepsilon} \\ \mathbf{y} \end{bmatrix} \sim MVN \left[\mathbf{0}, \begin{bmatrix} \mathbf{G}\sigma_u^2 & \mathbf{0} & \mathbf{G}\sigma_u^2 \\ \mathbf{0} & \mathbf{I}\sigma_\varepsilon^2 & \mathbf{I}\sigma_\varepsilon^2 \\ \mathbf{G}\sigma_u^2 & \mathbf{I}\sigma_\varepsilon^2 & \mathbf{G}\sigma_u^2 + \mathbf{I}\sigma_\varepsilon^2 \end{bmatrix} \right] \quad (S.4)$$

38 **Best Linear Unbiased Predictor of Genomic Values.** The conditional expectation (CE) is the best
 39 predictor in the mean-squared error sense, and in the model defined by (S.4) the CE of \mathbf{u} given \mathbf{y} is:

$$40 \quad \begin{aligned} E[\mathbf{u}|\mathbf{y}] &= Cov[\mathbf{u}, \mathbf{y}'] Var[\mathbf{y}]^{-1} \mathbf{y} \\ &= \sigma_u^2 \mathbf{G} [\mathbf{G}\sigma_u^2 + \mathbf{I}\sigma_\varepsilon^2]^{-1} \mathbf{y} \\ &= \mathbf{G} [\mathbf{G} + \mathbf{I}\lambda]^{-1} \mathbf{y} \\ &= \mathbf{G}\mathbf{T}\mathbf{y} \\ &= \mathbf{G}\tilde{\mathbf{y}} \end{aligned} \quad (S.5a)$$

41 where $\lambda = \sigma_\varepsilon^2 \sigma_u^{-2}$, $\mathbf{T} = [\mathbf{G} + \mathbf{I}\lambda]^{-1}$ is a matrix proportional to the inverse of the phenotypic variance-
 42 covariance matrix of phenotypes and $\tilde{\mathbf{y}} = \mathbf{T}\mathbf{y}$ is a vector of ‘smoothed’ phenotypes obtained by pre-
 43 multiplying \mathbf{y} with \mathbf{T} . The expected value (over possible realizations of \mathbf{u} and $\boldsymbol{\varepsilon}$) of the CE formula of

44 (S.5a) equal's the prior expectation of \mathbf{u} , $E\{E[\mathbf{u}|\mathbf{y}]\} = E[\mathbf{GTy}] = \mathbf{GTE}[\mathbf{y}] = \mathbf{0}$; therefore predictions
 45 derived using (S.5a) are unbiased with respect to the prior expectation. Finally, the predictor in (S.5a) is
 46 linear on \mathbf{y} ; therefore (S.5a) is the BLUP of \mathbf{u} .

47 The equation corresponding to the i^{th} entry in the CE vector of (S.5) is given by

$$48 \begin{aligned} E[u_i|\mathbf{y}] &= \sum_{j=1}^n G_{ij} \sum_{k=1}^n T_{jk} y_k \\ &= \sum_{j=1}^n G_{ij} \tilde{y}_j \end{aligned} \quad (\text{S.5b})$$

49 where $\tilde{y}_j = \sum_{k=1}^n T_{jk} y_k$.

50 **Best Linear unbiased Prediction of Marker Effects.** In the model above, marker effects and
 51 phenotypes follow the following MVN density:

$$52 \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{y} \end{bmatrix} \sim MVN \left[\mathbf{0}, \begin{pmatrix} \mathbf{I}\sigma_\beta^2 & \mathbf{W}'\sigma_\beta^2 \\ \mathbf{W}\sigma_\beta^2 & \mathbf{WW}'\sigma_\beta^2 + \mathbf{I}\sigma_\varepsilon^2 \end{pmatrix} \right]$$

53 Therefore, the expected value of marker effects given phenotypes is:

$$54 \begin{aligned} E[\boldsymbol{\beta}|\mathbf{y}] &= Cov[\boldsymbol{\beta}, \mathbf{y}'] Var[\mathbf{y}]^{-1} \mathbf{y} \\ &= \sigma_\beta^2 \mathbf{W}' [\mathbf{WW}'\sigma_\beta^2 + \mathbf{I}\sigma_\varepsilon^2]^{-1} \mathbf{y} \\ &= \frac{1}{p} \sigma_u^2 \mathbf{W}' [\mathbf{G}\sigma_u^2 + \mathbf{I}\sigma_\varepsilon^2]^{-1} \mathbf{y} \\ &= \frac{1}{p} \mathbf{W}' [\mathbf{G} + \mathbf{I}\lambda]^{-1} \mathbf{y} \\ &= \frac{1}{p} \mathbf{W}' \mathbf{T} \mathbf{y} = \frac{1}{p} \mathbf{W}' \tilde{\mathbf{y}} \end{aligned} \quad (\text{S.6})$$

55 Above, $\sigma_\beta^2 = \sigma_u^2 / p$, $\mathbf{G} = p^{-1} \mathbf{WW}'$ and $\mathbf{T} = [\mathbf{G} + \mathbf{I}\lambda]^{-1}$. It can be shown that predictions of marker
 56 effects given by (S.6) are equivalent to the Ridge Regression [3] estimates, $\hat{\boldsymbol{\beta}} = [\mathbf{WW}' + \mathbf{I}p\lambda]^{-1} \mathbf{W}' \mathbf{y}$.

57 **Note.** In the derivation of the BLUPs in (S.5) and (S.6) we have assumed that variance
 58 components are known. When variance components are unknown these can be estimated from data using
 59 various methods (e.g., Maximum Likelihood, Restricted Maximum Likelihood or Bayesian Methods, see
 60 Section 2 of this document). However, when variance components are estimated from the data, the CE

61 function does not have a closed form and predictions are not linear functions of the data; therefore,
 62 predictions derived when variance components are estimated from the data are not strictly BLUP
 63 anymore.

64 **Predictions of Yet-to-be Observed Phenotypes.** Consider a partition of the vectors in eq. (S.3)
 65 into two disjoint sets $\{\mathbf{y}_1, \mathbf{u}_1, \boldsymbol{\varepsilon}_1\}$ and $\{\mathbf{y}_2, \mathbf{u}_2, \boldsymbol{\varepsilon}_2\}$, pertaining to the training (TRN) and testing (TST) data
 66 sets, respectively, so that:

$$67 \quad \begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{u}_1 \\ \mathbf{u}_2 \end{bmatrix} + \begin{bmatrix} \boldsymbol{\varepsilon}_1 \\ \boldsymbol{\varepsilon}_2 \end{bmatrix} \quad (\text{S.7})$$

68 The joint density of phenotypes in TRN and TST data sets is then given by the following MVN density:

$$69 \quad \begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{bmatrix} \sim MVN \left[\mathbf{0}, \begin{pmatrix} \mathbf{G}_{11}\sigma_u^2 + \mathbf{I}_1\sigma_\varepsilon^2 & \mathbf{G}_{12}\sigma_u^2 \\ \mathbf{G}_{21}\sigma_u^2 & \mathbf{G}_{22}\sigma_u^2 + \mathbf{I}_2\sigma_\varepsilon^2 \end{pmatrix} \right] \quad (\text{S.8})$$

70

71 where \mathbf{G}_{11} and \mathbf{G}_{22} represent matrices describing genomic relationships among individuals of the TRN
 72 and TST data sets, respectively, and $\mathbf{G}_{21} = \mathbf{G}'_{12}$ is a matrix describing genomic relationships between
 73 individuals in the TST and those in the TRN data set. From equation (S.8) the expected value of
 74 phenotypes in the TST data set given the phenotypes in the TRN data set is,

$$75 \quad \begin{aligned} E[\mathbf{y}_2|\mathbf{y}_1] &= Cov[\mathbf{y}_2, \mathbf{y}_1] Var[\mathbf{y}_1]^{-1} \mathbf{y}_1 \\ &= \sigma_u^2 \mathbf{G}_{21} [\sigma_u^2 \mathbf{G}_{11} + \sigma_\varepsilon^2 \mathbf{I}_1]^{-1} \mathbf{y}_1 \\ &= \mathbf{G}_{21} [\mathbf{G}_{11} + \mathbf{I}_1 \lambda]^{-1} \mathbf{y}_1 \\ &= \mathbf{G}_{21} \mathbf{T} \mathbf{y}_1 \\ &= \mathbf{G}_{21} \tilde{\mathbf{y}}_1 \end{aligned} \quad (\text{S.9a})$$

76 Above, $\mathbf{T} = [\mathbf{G}_{11} + \mathbf{I}_1 \lambda]^{-1}$ is a matrix proportional to the inverse of the (co)variance matrix of phenotypes
 77 in the TRN data set. A scalar version of (S.9a) is given by the following linear score (hereinafter, we
 78 assume that the TRN data set includes n individuals and, without loss of generality, we present prediction
 79 equations and other expression pertaining genetic values or phenotypes in the TST data set using sub-
 80 index $n+1$ to stress that such expression pertains to observations not included in the TRN data set):

81

$$\begin{aligned}
E[y_{n+1}|\mathbf{y}_1] &= \sum_{i=1}^n G_{n+1,i} \sum_{j=1}^n T_{ij} y_j \\
&= \sum_{i=1}^n G_{n+1,i} \tilde{y}_i
\end{aligned} \tag{S.9b}$$

82 Therefore, predictions in G-BLUP are simply weighted sums of (pre-smoothed) phenotypes of the
83 individuals in the TRN data set, with weights given by the realized genomic relationships between
84 individuals in TRN and TST data sets $\{G_{n+1,i}\}_{i=1}^n$.

85

86 2. Measures of Prediction Accuracy in G-BLUP

87

88 The variance of prediction errors, or prediction error variance (PEV), is commonly used in prediction
89 problems to assess the predictive power of an entertained model. In case of prediction of yet-to-be-
90 observed phenotypes, variance-covariance matrices of prediction errors of genetic values and of
91 phenotypes are given by

$$92 \quad \text{Var}(\mathbf{u}_2 - E[\mathbf{u}_2|\mathbf{y}_1]) = \text{Var}(\mathbf{u}_2) + \text{Var}(E[\mathbf{u}_2|\mathbf{y}_1]) - 2\text{Cov}(\mathbf{u}_2, E[\mathbf{u}_2|\mathbf{y}_1])$$

93

and

$$94 \quad \text{Var}(\mathbf{y}_2 - E[\mathbf{y}_2|\mathbf{y}_1]) = \text{Var}(\mathbf{y}_2) + \text{Var}(E[\mathbf{y}_2|\mathbf{y}_1]) - 2\text{Cov}(\mathbf{y}_2, E[\mathbf{y}_2|\mathbf{y}_1])$$

95

96 , respectively. The next sub-section presents expressions for prediction error variances and co-variances
under the assumptions of the model.

97

98

99 **2.1. Case 1: when genotypes at causal loci are known (i.e., perfect LD between markers and QTL)**
100

101 We begin by assuming that genomic relationship matrices are computed using realized genotypes at
102 causal loci. In this case, the model holds, and closed-form formulas for PEV and R^2 can be derived. The
103 results that follow use (S.8) as starting point and standard properties of the multivariate normal density. A
104 list of useful results that will be used later is given in Table S1.

105

106 **Table S1.** Useful Results: variances, (co)variances and prediction error variances of genetic values and of
107 phenotypes of individuals in the testing data set.

	Genetic Values	Phenotypes
Marginal Variance	$Var(\mathbf{u}_2) = \mathbf{G}_{22}\sigma_u^2$ (S10.a)	$Var(\mathbf{y}_2) = [\mathbf{G}_{22}\sigma_u^2 + \mathbf{I}_2\sigma_\varepsilon^2]$ (S10.b)
Conditional Variances	$Var(\mathbf{u}_2 \mathbf{y}_1) = \sigma_u^2\{\mathbf{G}_{22} - \mathbf{G}_{21}\mathbf{T}\mathbf{G}_{12}\}$ (S10.c)	$Var(\mathbf{y}_2 \mathbf{y}_1) = Var(\mathbf{u}_2 \mathbf{y}_1) + \mathbf{I}_2\sigma_\varepsilon^2$ (S10.d)
Variance of Cond. Expectation	$Var(E[\mathbf{y}_2 \mathbf{y}_1]) = Var(E[\mathbf{u}_2 \mathbf{y}_1]) = \sigma_u^2\mathbf{G}_{21}\mathbf{T}\mathbf{G}_{12}$ (S10.e)	
Cov. Predictions and realized values	$Cov(\mathbf{u}_2, E[\mathbf{u}_2 \mathbf{y}_1]) = Cov(\mathbf{y}_2, E[\mathbf{y}_2 \mathbf{y}_1]) = \sigma_u^2\mathbf{G}_{21}\mathbf{T}\mathbf{G}_{12}$ (S10.e)	
Prediction Error Variances	$Var(\mathbf{u}_2 - E[\mathbf{u}_2 \mathbf{y}_1]) = \sigma_u^2[\mathbf{G}_{22} - \mathbf{G}_{21}\mathbf{T}\mathbf{G}_{12}]$ (S.10c)	$Var(\mathbf{y}_2 - E[\mathbf{y}_2 \mathbf{y}_1]) = \sigma_u^2[\mathbf{G}_{22} - \mathbf{G}_{21}\mathbf{T}\mathbf{G}_{12}] + \mathbf{I}_2\sigma_\varepsilon^2$ (S.10d)

108

109 The first element on the right hand side of (S.10d), is the variance-covariance matrix of prediction errors
110 of genetic values, $Var(\mathbf{u}_2 - E[\mathbf{u}_2|\mathbf{y}_1]) = \sigma_u^2[\mathbf{G}_{22} - \mathbf{G}_{21}\mathbf{T}\mathbf{G}_{12}]$; this is simply the prior variance-
111 covariance matrix of genetic values, $\sigma_u^2\mathbf{G}_{22}$ minus a quadratic form, given by $\sigma_u^2\mathbf{G}_{21}\mathbf{T}\mathbf{G}_{12}$, that
112 quantifies reduction in uncertainty about genetic values of individuals in the TST data set attained by
113 observing phenotypes of the individuals in the TRN data set. If individuals in TRN and TST data sets, are
114 independent, i.e., when $\mathbf{G}_{21} = \mathbf{0}$, the (co)variance of prediction errors of genetic values equals the prior
115 variance (i.e., there is no statistical learning). The second term in the right-hand-side of (S.10d), $\mathbf{I}_2\sigma_\varepsilon^2$,
116 represents uncertainty about future phenotypes due to error terms; since these are uncorrelated with the
117 phenotypes in the TRN data set there is no learning about this component.

118 The diagonal elements of expressions S.10c and S.10d give the PEV of individual genetic values
 119 and of individual phenotypes, respectively,

$$120 \quad \text{Var}(u_{n+1} - E[u_{n+1}|\mathbf{y}_1]) = \sigma_u^2 \left[G_{n+1,n+1} - \sum_{i=1}^n \sum_{j=1}^n G_{n+1,i} G_{n+1,j} T_{ij} \right] \quad (\text{S.10e})$$

$$121 \quad \text{Var}(y_{n+1} - E[y_{n+1}|\mathbf{y}_1]) = \sigma_u^2 \left[G_{n+1,n+1} - \sum_{i=1}^n \sum_{j=1}^n G_{n+1,i} G_{n+1,j} T_{ij} \right] + \sigma_\varepsilon^2 \quad (\text{S.10f})$$

122 Since \mathbf{T} is positive definite $\sum_{i=1}^n \sum_{j=1}^n G_{n+1,i} G_{n+1,j} T_{ij} \geq 0$ with equality holding only when all the genomic
 123 relationships between individual $n+1$ and individuals in the TRN data set are equal to zero.

124 **R-squared.** The proportional reduction of variance of genetic values accounted for by predictions
 125 is given by the following R-squared (genetic values):

$$\begin{aligned} R_{n+1,u}^2 &= 1 - \frac{\text{Var}(u_{n+1} - E[u_{n+1}|\mathbf{y}_1])}{\text{Var}(u_{n+1})} \\ &= 1 - \frac{\sigma_u^2 \left[G_{n+1,n+1} - \sum_{i=1}^n \sum_{j=1}^n G_{n+1,i} G_{n+1,j} T_{ij} \right]}{\sigma_u^2 G_{n+1,n+1}} \\ &= \frac{\sum_{i=1}^n \sum_{j=1}^n G_{n+1,i} G_{n+1,j} T_{ij}}{G_{n+1,n+1}} \end{aligned} \quad (\text{S.11})$$

127 A similar R-squared for prediction of phenotypes is:

$$\begin{aligned} R_{n+1,y}^2 &= 1 - \frac{\text{Var}(y_{n+1} - E[y_{n+1}|\mathbf{y}_1])}{\text{Var}(y_{n+1})} \\ &= 1 - \frac{\sigma_u^2 \left[G_{n+1,n+1} - \sum_{i=1}^n \sum_{j=1}^n G_{n+1,i} G_{n+1,j} T_{ij} \right] + \sigma_\varepsilon^2}{\sigma_u^2 G_{n+1,n+1} + \sigma_\varepsilon^2} \\ &= \frac{\sigma_u^2 \sum_{i=1}^n \sum_{j=1}^n G_{n+1,i} G_{n+1,j} T_{ij}}{\sigma_u^2 G_{n+1,n+1} + \sigma_\varepsilon^2} \end{aligned}$$

129 The ratio $R_{n+1,y}^2 / R_{n+1,u}^2 = \frac{\sigma_u^2 G_{n+1,n+1}}{\sigma_u^2 G_{n+1,n+1} + \sigma_\varepsilon^2} = h_{n+1}^2$ where h_{n+1}^2 is the heritability of the trait for individual

130 $n+1$ (note that, when $G_{n+1,n+1} = 1$, h_{n+1}^2 equals the heritability of the trait, $h^2 = \frac{\sigma_u^2}{\sigma_u^2 + \sigma_\varepsilon^2}$); therefore,

131
$$R_{n+1,y}^2 = h_{n+1}^2 R_{n+1,u}^2 = h_{n+1}^2 \frac{\sum_{i=1}^n \sum_{j=1}^n G_{n+1,i} G_{n+1,j} T_{ij}}{G_{n+1,n+1}} \quad (\text{S.12})$$

132 Suppose now that all the off-diagonal elements of \mathbf{G}_{11} (i.e., all the G_{ij} such that
 133 $i, j \in (1, \dots, n); i \neq j$) are zero. In this case \mathbf{T} is diagonal and the R-squared formula reduces to

134
$$R_{n+1,y}^2 = h_{n+1}^2 \frac{\sum_{i=1}^n G_{n+1,i}^2}{G_{n+1,n+1}}.$$
 When data involve nominally unrelated individuals the off-diagonal elements

135 of \mathbf{G}_{11} are expected to be small, but not exactly equal to zero. Relative to the case where the TRN data set
 136 comprise independent TRN phenotypes, the use of correlated TRN phenotypes yields smaller R-squared.
 137 This happens because, other things being equal, the amount of information provided by the TRN data set

138 comprises genetically independent individuals. Therefore, the index $\frac{h_{n+1}^2}{G_{n+1,n+1}} \sum_{i=1}^n G_{n+1,i}^2$ acts as an upper-

139 bound to the R-squared; specifically,

140
$$\text{Case 1 (perfect LD between markers and QTL): } R_{n+1,y}^2 \leq \frac{h_{n+1}^2}{G_{n+1,n+1}} \sum_{i=1}^n G_{n+1,i}^2. \quad (\text{S.13})$$

141

142 **2.2. Case 2: when genomic relationships are computed using markers (i.e., imperfect LD between**
 143 **markers and QTL)**

144 In practice, genomic relationships are computed from marker genotypes that are in imperfect LD
 145 with genotypes at causal loci. Further, since the patterns of realized genomic relationships at different sets
 146 of loci (e.g., markers and causal loci) vary across the genome [4], realized genomic relationships at
 147 markers ($\overline{G}_{n+1,i}$) should be regarded as proxies for the realized genomic relationships at causal loci (
 148 $G_{n+1,i}$).

149 When marker-derived genomic relationships are used in place of realized genomic relationships
 150 at causal loci, the assumptions of model (S.8) do not hold. Therefore it is not possible to derive closed-
 151 form expressions for prediction R-squared. This problem can be circumvented by deriving a closed form
 152 expression for an upper bound to prediction R-squared. To arrive at this closed form we assume that
 153 genomic relationships realized at causal loci among pairs of individuals in the TRN data set are known.
 154 Essentially, this has the effect of treating matrix \mathbf{T} as a known constant. Therefore, we consider only the

155 impacts of imperfect LD between markers and QTL that occurs through misspecification of genomic
 156 relationships between individuals in the TST and those in the TRN data set. Predictions in G-BLUP are
 157 given by $\sum_{i=1}^n G_{n+1,i} \tilde{y}_i$. Because of the assumption that genomic relationships at causal loci among
 158 individuals in the TRN data set are known, inferences about the \tilde{y}_i 's, the entries of the vector
 159 $\tilde{\mathbf{y}} = [\mathbf{G} + \mathbf{I}\lambda]^{-1} \mathbf{y}$, are not affected by imperfect LD between markers and QTL.

160 Assume that the realized genomic relationships at markers between an individual in the TST data
 161 set ($n+1$) and all individuals in the TRN data set, $\{\bar{G}_{n+1,1}, \dots, \bar{G}_{n+1,n}\}$, can be described using a linear
 162 regression on genomic relationships realized at causal loci, $\{G_{n+1,1}, \dots, G_{n+1,n}\}$, that is

$$163 \quad \bar{G}_{n+1,i} = b_{n+1} G_{n+1,i} + \xi_{n+1,i}, \quad (i=1, \dots, n)$$

164 where b_{n+1} represents the regression of marker-derived genomic relationships $\{\bar{G}_{n+1,1}, \dots, \bar{G}_{n+1,n}\}$ on
 165 realized genomic relationships at causal loci $\{G_{n+1,1}, \dots, G_{n+1,n}\}$ and $\xi_{n+1,i}$ represents an error term which
 166 accounts for differences in realized genomic relationships at markers and QTL. A similar approach was
 167 used before by Yang and co-authors [5]. However, the regression used by these author applied to all the
 168 entries of the relationship matrix, including diagonal and off-diagonal terms. In our case, we focus on
 169 quantifying the effects of imperfect LD that occur through miss-specification of TRN-TST relationships;
 170 therefore the regression is based on the off-diagonal terms only.

171 Using $\bar{G}_{n+1,i} = b_{n+1} G_{n+1,i} + \xi_{n+1,i}$ in (S.9b) we get:

$$172 \quad \begin{aligned} \bar{E}[u_{n+1} | \mathbf{y}_1] &= \sum_{i=1}^n \bar{G}_{n+1,i} \sum_{j=1}^n T_{ij} y_j \\ &= \sum_{i=1}^n \bar{G}_{n+1,i} \tilde{y}_i \\ &= \sum_{i=1}^n (b_{n+1} G_{n+1,i} + \xi_{n+1,i}) \tilde{y}_i \end{aligned}$$

173 The term $\xi_{n+1,i}$ is, by construction, orthogonal to the realized genomic relationships at causal loci;
 174 therefore, for large n it is safe to assume that $\sum_{i=1}^n \xi_{n+1,i} \tilde{y}_i$ approaches zero. Using $\sum_{i=1}^n \xi_{n+1,i} \tilde{y}_i = 0$ in
 175 the above expression, $\bar{E}[u_{n+1} | \mathbf{y}_1] = b_{n+1} \sum_{i=1}^n G_{n+1,i} \tilde{y}_i = b_{n+1} E[u_{n+1} | \mathbf{y}_1]$. Therefore, the variance of
 176 prediction errors is now given by

$$\begin{aligned}
177 \quad \text{Var}(u_{n+1} - \bar{E}[u_{n+1}|\mathbf{y}_1]) &= \text{Var}(u_{n+1} - b_{n+1}E[u_{n+1}|\mathbf{y}_1]) \\
&= \text{Var}(u_{n+1}) + b_{n+1}^2 \text{Var}(E[u_{n+1}|\mathbf{y}_1]) - 2b_{n+1} \text{Cov}(u_{n+1}, E[u_{n+1}|\mathbf{y}_1])
\end{aligned}$$

178 It can be shown (see expression (S10.e) in Table S1) that the variance of the conditional expectation,
179 $\text{Var}(E[u_{n+1}|\mathbf{y}_1])$, equals the covariance between the true genetic values and the BLUP,
180 $\text{Cov}(u_{n+1}, E[u_{n+1}|\mathbf{y}_1])$. Specifically, from (S10.e) we have:

$$181 \quad \text{Var}(E[u_{n+1}|\mathbf{y}_1]) = \text{Var}[u_{n+1}|\mathbf{y}_1] = \text{Cov}(u_{n+1}, E[\mathbf{u}_2|\mathbf{y}_1]) = \sigma_u^2 \mathbf{g}'_{n+1} \mathbf{T} \mathbf{g}_{n+1} = \sigma_u^2 \sum_{i=1}^n \sum_{j=1}^n G_{n+1,i} G_{n+1,j} T_{ij}$$

182 , where $\mathbf{g}_{n+1} = \{G_{n+1,1}, \dots, G_{n+1,n}\}^{\text{c}}$ is a vector containing the genomic relationships realized at causal loci
183 between individual $n+1$ and every individual in the TRN data set. This is simply a scalar version of
184 (S10.e). Therefore,

$$\begin{aligned}
\text{Var}(u_{n+1} - \bar{E}[u_{n+1}|\mathbf{y}_1]) &= \text{Var}(u_{n+1}) + b_{n+1}^2 \text{Var}(E[u_{n+1}|\mathbf{y}_1]) - 2b_{n+1} \text{Cov}(u_{n+1}, E[u_{n+1}|\mathbf{y}_1]) \\
185 \quad &= \sigma_u^2 \left\{ G_{n+1,n+1} + (b_{n+1}^2 - 2b_{n+1}) \sum_{i=1}^n \sum_{j=1}^n G_{n+1,i} G_{n+1,j} T_{ij} \right\} \\
&= \sigma_u^2 \left\{ G_{n+1,n+1} - (2b_{n+1} - b_{n+1}^2) \sum_{i=1}^n \sum_{j=1}^n G_{n+1,i} G_{n+1,j} T_{ij} \right\}
\end{aligned}$$

186 Using $(2b_{n+1} - b_{n+1}^2) = 1 - (1 - b_{n+1})^2$ the prediction R^2 becomes

$$\begin{aligned}
\bar{R}_{n+1,u}^2 &= 1 - \frac{\text{Var}(u_{n+1} - \bar{E}[u_{n+1}|\mathbf{y}_1])}{\text{Var}(u_{n+1})} \\
187 \quad &= 1 - \frac{\sigma_u^2 \left\{ G_{n+1,n+1} - [1 - (1 - b_{n+1})^2] \sum_{i=1}^n \sum_{j=1}^n G_{n+1,i} G_{n+1,j} T_{ij} \right\}}{\sigma_u^2 G_{n+1,n+1}} \\
&= [1 - (1 - b_{n+1})^2] \frac{\sum_{i=1}^n \sum_{j=1}^n G_{n+1,i} G_{n+1,j} T_{ij}}{G_{n+1,n+1}} \\
&= [1 - (1 - b_{n+1})^2] R_{n+1,u}^2
\end{aligned}$$

188 And for prediction of phenotypes we have:

$$189 \quad \bar{R}_{n+1,y}^2 = h_{n+1}^2 [1 - (1 - b_{n+1})^2] R_{n+1,u}^2 = [1 - (1 - b_{n+1})^2] R_{n+1,y}^2$$

$$190 \quad \text{and, } \frac{R_{n+1,y}^2 - \bar{R}_{n+1,y}^2}{R_{n+1,y}^2} = (1 - b_{n+1})^2.$$

191 Therefore, the coefficient $\tau_{n+1} = (1 - b_{n+1})^2$, can be regarded as a minimum shrinkage factor on R-
 192 squared due to use of marker-derived genomic relationships, as opposed to those realized at causal loci.
 193 This is a minimum shrinkage factor because in its derivation we have assumed that genomic relationships
 194 at causal loci were known for individuals at the TRN data set. Hence, we have

195 **Case 2 (imperfect LD):** $\bar{R}_{n+1,y}^2 \leq h_{n+1}^2 [1 - (1 - b_{n+1})^2] R_{n+1,u}^2$ (S.14)

196 Under perfect LD and with infinite sample size, $R_{n+1,u}^2$ reaches one. Therefore, the term
 197 $h_{n+1}^2 [1 - (1 - b_{n+1})^2]$ can be interpreted as an ‘optimistic’ upper bound on R-squared of predictions of
 198 phenotypes that be attained using G-BLUP type methods.

199 The term $\tau_{n+1} = (1 - b_{n+1})^2$ plays a central role in upper-bounds on R-squared. This coefficient has a
 200 maximum at $b_{n+1} = 1$, a case that would occur if markers and causal loci and in perfect LD. Values of
 201 b_{n+1} smaller or larger than one yield $[1 - (1 - b_{n+1})^2] < 1$, reducing the maximum R-squared that can be
 202 attained.

203 In practice, the set of causal loci is unknown; however, an approximation to b_{n+1} can be obtained
 204 by computing realized genomic relationships, $\{\bar{G}_{n+1,1}, \dots, \bar{G}_{n+1,n}\}$ and $\{G_{n+1,l}, \dots, G_{n+1,n}\}$, at disjoint sets of
 205 loci and regressing $\{\bar{G}_{n+1,l}, \dots, \bar{G}_{n+1,n}\}$ on $\{G_{n+1,l}, \dots, G_{n+1,n}\}$. We present estimates of these coefficients in
 206 the article and discuss the impact of linkage on this regression by comparing estimates of these
 207 regressions derived using data with related and unrelated individuals.

208

209 **1. Bayesian Implementation**

210

211 The model described by expression S.4 was implemented with two modifications: (a) the model was
 212 extended by inclusion of an intercept, and (b) variance parameters were treated as unknown and were
 213 estimated from the training data sets. The model was implemented in a Bayesian setting; this requires
 214 assigning prior distributions to the unknown intercept and variance components. The intercept was
 215 assigned a flat prior and variance components were assigned independent Scaled-Inverse Chi-square
 216 densities. Therefore, the joint posterior density of the unknowns was

217
$$p(\mu, \mathbf{u}, \sigma_u^2, \sigma_\varepsilon^2 | \mathbf{y}, df, S_\varepsilon, S_u) \propto p(\mathbf{y} | \mathbf{u}, \sigma_\varepsilon^2) p(\mathbf{u} | \sigma_u^2) p(\sigma_u^2, \sigma_\varepsilon^2)$$

218
$$\propto \text{MVN}(\mathbf{y} | \mathbf{1}\mu + \mathbf{u}, \mathbf{I}\sigma_\varepsilon^2) \text{MVN}(\mathbf{u} | \mathbf{0}, \mathbf{G}\sigma_u^2) \chi^{-2}(\sigma_\varepsilon^2 | df, S_\varepsilon) \chi^{-2}(\sigma_u^2 | df, S_u) \quad (\text{S.15})$$

218 where, $\chi^{-2}(\sigma^2 | df, S)$ denotes an Inverted-Scaled Chi-Squared density for the random variable σ^2 and
 219 with degree of freedom and scale parameters df and S , respectively. The degree of freedom parameter
 220 was set equal to 5 and the scale was chosen so that the prior expected values of the variance parameters
 221 equals the phenotypic variance (1 in the simulated data and approximately 40 in the real data) times 0.8
 222 for the genomic variance and times 0.2 in case of the residual variance.

223 Samples from the posterior density were obtained with a Gibbs sampler that uses an orthogonal
 224 representation of the model (i.e., a random regression on the eigenvectors of \mathbf{G}). The algorithm is fully
 225 described in [6,7]. A total of 18,000 samples per run were drawn. Of these, the first 3,000 were discarded
 226 as burn-in. The remaining 15,000 samples were thinned with a thinning interval of 3, yielding total of
 227 5,000 samples to compute Monte Carlo estimates of the unknowns. The software used to carry out
 228 analyses is available upon request to the corresponding author.

229

230 Literature Cited

231 1. Henderson CR (1975) Best linear unbiased estimation and prediction under a selection model.
 232 Biometrics 31: 423–447.

233 2. Pszczola M, Strabel T, Mulder HA, Calus MPL (2012) Reliability of direct genomic values for
 234 animals with different relationships within and to the reference population. Journal of dairy science
 235 95: 389–400.

236 3. Hoerl AE, Kennard RW (1970) Ridge regression: Biased estimation for nonorthogonal problems.
 237 Technometrics 12: 55–67.

238 4. Hill WG, Weir BS (2011) Variation in actual relationship as a consequence of Mendelian sampling
 239 and linkage. Genetics Research 93: 47–64. doi:10.1017/S0016672310000480.

240 5. Yang J, Benyamin B, McEvoy BP, Gordon S, Henders AK, et al. (2010) Common SNPs explain a
 241 large proportion of the heritability for human height. Nature genetics 42: 565–569.

242 6. de los Campos G, Gianola D, Rosa GJM, Weigel KA, Crossa J (2010) Semi-parametric genomic-
 243 enabled prediction of genetic values using reproducing kernel Hilbert spaces methods. Genetics
 244 Research 92: 295–308.

245 7. Janss L, de los Campos G, Sheehan N, Sorensen DA (2012) Inferences from Genomic Models in
 246 Stratified Populations. GENETICS: 693–704. doi:10.1534/genetics.112.141143.

- 247 1. Henderson CR (1975) Best linear unbiased estimation and prediction under a selection model.
248 Biometrics 31: 423–447.
- 249 2. Pszczola M, Strabel T, Mulder HA, Calus MPL (2012) Reliability of direct genomic values for
250 animals with different relationships within and to the reference population. Journal of dairy science
251 95: 389–400.
- 252 3. Hoerl AE, Kennard RW (1970) Ridge regression: Biased estimation for nonorthogonal problems.
253 Technometrics 12: 55–67.
- 254 4. Hill WG, Weir BS (2011) Variation in actual relationship as a consequence of Mendelian sampling
255 and linkage. Genetics Research 93: 47–64. doi:10.1017/S0016672310000480.
- 256 5. Yang J, Benyamin B, McEvoy BP, Gordon S, Henders AK, et al. (2010) Common SNPs explain a
257 large proportion of the heritability for human height. Nature genetics 42: 565–569.
- 258 6. De los Campos G, Gianola D, Rosa GJM, Weigel KA, Crossa J (2010) Semi-parametric genomic-
259 enabled prediction of genetic values using reproducing kernel Hilbert spaces methods. Genetics
260 Research 92: 295–308.
- 261 7. Janss L, de los Campos G, Sheehan N, Sorensen DA (2012) Inferences from Genomic Models in
262 Stratifi_ed Populations. GENETICS: 693–704. doi:10.1534/genetics.112.141143.

263