



This is the pre-peer reviewed version of the following article:

Gowane, G., Lee, S., Clark, S., Moghaddar, N., Al-Mamun, H., & Werf, J. (2019). Effect of selection and selective genotyping for creation of reference on bias and accuracy of genomic prediction. *Journal Of Animal Breeding And Genetics*, 136(5), 390-407. doi: 10.1111/jbg.12420

which has been published in final form at <https://doi.org/10.1111/jbg.12420>. This article may be used for non-commercial purposes in accordance with Wiley Terms and Conditions for Use of Self-Archived Versions."

This work is licensed under an [Attribution-NonCommercial-NoDerivatives 4.0 International licence \(CC BY-NC-ND 4.0\)](#)

1 **Effect of selection on bias and accuracy in genomic prediction of breeding values**

2 *G. R. Gowane¹, Sang Hong Lee³, Sam Clark², Nasir Moghaddar², Hawlader A Al- Mamun⁴ & Julius H. J.*
3 *van der Werf²*

4 ¹*Animal Genetics & Breeding Division, ICAR-Central Sheep & Wool research Institute, Avikanagar*
5 *304501 Rajasthan, India*

6 gopalgowane@gmail.com (Corresponding Author)

7 ²*School of Environmental and Rural Sciences, University of New England, Armidale 2351 NSW, Australia*

8 ³*Centre for Population Health Research, School of Health Sciences and Sansom Institute of Health*
9 *Research, University of South Australia, Adelaide, Australia*

10 ⁴*Statistical Genetics and Genomics, DATA61/CSIRO, GPO Box 1700, Canberra, ACT 2601, Australia*

11

12 **Abstract**

13 Reference populations for genomic selection (GS) usually involve highly selected individuals,
14 which may result in biased prediction of estimated genomic breeding values (GEBV). In the present
15 study, bias and accuracy of GEBV were explored for various genetic models and prediction methods
16 when using selected individuals for a reference. Data were simulated for an animal breeding program to
17 compare Best Linear Unbiased Prediction of breeding values using pedigree based relationships
18 (PBLUP), genomic relationships for genotyped animals only (GBLUP) and a Single Step approach
19 (SSGBLUP), where information on genotyped individuals was used to infer a matrix **H** with relationships
20 among all available genotyped and non-genotyped individuals that were linked through pedigree. In
21 SSGBLUP, various weights ($\alpha=0.95, 0.80, 0.50$) for the genomic relationship matrix (**G**) relative to the
22 numerator relationship matrix (**A**) were applied to construct **H** and in another version (SSGBLUP_F),
23 inbreeding was accounted for while computing \mathbf{A}^{-1} . With GBLUP, accuracy of GEBV prediction
24 increased linearly with an increase in the number of animals selected in reference. For the scenario with
25 no-selection and random mating (RR) prediction was unbiased. For GBLUP, lower accuracy and bias

1 observed in the scenarios with selection and random mating (SR) or selection and positive assortative
2 mating (SA), in which prediction bias increased when a smaller and highly selected proportion
3 genotyped. Bias disappeared when all individuals were genotyped. SSGBLUP_F showed higher accuracy
4 compared to GBLUP and bias of prediction was negligible even with selective genotyping. However,
5 PBLUP and SSGBLUP showed bias in SA owing to not fully accounting for allele frequency changes
6 because of selection of quantitative trait loci (QTL) with larger effects and also due to high inbreeding
7 rate. In genetic models with fewer QTL but each with larger effect, predictions were less accurate and
8 more biased for selection scenarios. Results suggest that prediction accuracy and bias is affected by the
9 genetic architecture of the trait. Selective genotyping lead to significant bias in GEBV prediction.
10 SSGBLUP with appropriate scaling of **A** and **G** matrices can provide accurate and less biased prediction
11 but scaling requires careful consideration in populations under selection and with high levels of
12 inbreeding.

13 **Key Words: Genomic Selection, GBLUP, Prediction Bias, Single Step Evaluation**

14 **Introduction**

15 Best Linear Unbiased Prediction (BLUP) provides unbiased estimates of breeding values in populations
16 under selection, conditional on the inclusion of all information used in the selection decisions (Henderson
17 1975; Sorensen and Kennedy 1984). In the case of genomic prediction based on a selected reference
18 population, this condition might not be met, in particular when only genotyped individuals are evaluated
19 in genomic BLUP (GBLUP). A number of studies (VanRaden et al., 2009a, 2009b, Patry and Ducrocq
20 2011, Vitezica et al. 2011) reported decreased accuracy of genomic estimated breeding value (GEBV)
21 and increased bias due to selective genotyping of sires.

22 Single Step genomic BLUP (SSGBLUP) (Legarra et al. 2009; Christensen and Lund 2010) combines
23 genomic relationships from genotyped individuals with pedigree relationships with non-genotyped
24 individuals and this integration should allow information on unselected animals to be included, with all
25 relationships tracing back to a conceptual unselected base population. However, SSGBLUP requires the
26 genomic relationship matrix (**G**) and pedigree-based relationship matrix (**A**) to refer to the same base
27 population as otherwise new bias could be introduced. Some studies have discussed this issue and

1 proposed a slight modification in the SSGBLUP procedure (Forni et al. 2011, Vitezica et al. 2011,
2 Christensen et al. 2012). The modification suggested by Vitezica et al. (2011) involved adding the
3 difference between means of \mathbf{A}_{22} and \mathbf{G} to all elements of \mathbf{G} , such that the relationships of genotyped
4 individuals are rightly scaled in relation to the base population of all animals in the pedigree. In their
5 study, Vitezica et al. (2011) used a model with only 250 quantitative trait loci (QTL) affecting the trait
6 and 10 generations of assortative mating of sires and dam was used.

7 The present study aimed at exploring the effect of selection on genomic prediction for a wider range of
8 scenarios. Genetic evaluations based on pedigree BLUP (PBLUP), GBLUP and SSGBLUP were
9 compared. We investigated a number of factors affecting bias and accuracy of genetic evaluations,
10 including (1) the proportion of individuals selected to have genotype information; (2) the genetic
11 structure of the trait as determined by the heritability and the number of QTLs explaining the variance in
12 the trait; and (3) scenarios with and without selection and assortative mating.

13 **Material and Methods**

14 *Population and genotype simulation*

15 Data were simulated using QMSim (Sargolzaei & Schenkel, 2009). A historical population with effective
16 population size (N_e) of 100 was generated with 50 males and 50 females producing 2 progeny each by
17 random union of gametes in each of 95 generations and thereafter the number of progeny was gradually
18 expanded to 1000 offspring until the 100th generation. In the subsequent 10 generations (101-110), 50
19 males were mated to 500 females who produced 1000 progeny. The genomic structure consisted of 30
20 chromosomes of equal length (1 Morgan). Biallelic markers (60,000) were randomly distributed across
21 the genome with an equal frequency (0.5) in the first generation of the historical population. The
22 mutation rate of the markers and QTL was 2.5×10^{-8} per locus per generation (Hickey and Gorjanc 2012).
23 In order to make sure that data were correctly simulated for the historical population, we confirmed
24 whether population parameters like effective number of chromosome segments (M_e) in the simulated data
25 agreed with the expectation from theory given the value of N_e and the family structure in the population
26 (Lee et al. 2017). M_e was determined from the variation in genomic relationship among members of the
27 population. Four genetic models were simulated that differed in the number and distribution of QTL

1 effects; with 90, 990, 9,990 and 60,000 QTLs. The QTL allele effects were sampled from a Gamma
2 distribution with a shape and scale parameter of 0.4 and 1.0, respectively. Genotype effects at individual
3 QTL were aggregated to form true breeding values (TBV) and these were re-scaled to match the input
4 value for the additive genetic variance. For individuals in the 101st generation these breeding values were
5 normally distributed with mean 0 and variance h^2 . Residual effects on phenotype were independent and
6 normally distributed with mean 0 and variance $(1-h^2)$. Therefore, the mean and variance for the simulated
7 phenotypes were zero and one, respectively. Phenotypes were created for all individuals in the last 10
8 generations. Three different values for heritability ($h^2 = 0.1, 0.3$ and 0.5) were tested for each scenario.
9 Selection of sires was either random or based on estimated breeding values (EBV) obtained by PBLUP.
10 Positive assortative mating was also applied, giving rise to three different scenarios: I: Random selection
11 and random mating (RR), II: Selection on EBV and Random mating (SR), III: Selection on EBV and
12 assortative mating based on EBV (SA). Selection was with non-overlapping generations and in each
13 generation 10% of males were mated to all young females. Under each of these 3 scenarios, predictions
14 of EBV were obtained by each of the different methods and modelling parameters. Scenario PBLUP
15 involved BLUP prediction based on all phenotype and pedigree information from 9 generations to predict
16 EBV in the 10th generation. For each of the last 4 generations (6th to 9th generation), either 125 (25%), 250
17 (50%) or 500 (100%) males were selected for genotyping to form a reference population consisting of
18 500, 1000 or 2000 animals (scenario G500, G1000 or G2000). In scenario G9550, all the 9550 animals
19 from all 9 preceding generations were genotyped and used as a reference population. In scenario SS500,
20 SS1000 and SS2000, SSGBLUP was used for prediction of breeding values combining information from
21 pedigree (9550 individual's pedigree) and genomic relationships from either 500, 1000 or 2000
22 genotyped individuals in the reference, respectively. For each scenario, we conducted 25 replicates.

23 *Analysis and estimation of breeding value*

24 In total 1000 selection candidates in the 10th generation were used as a validation population to determine
25 the bias and accuracy of estimated breeding values.

26 PBLUP analysis was based on pedigree information from the 9 preceding generations that include 9550
27 animals. Alternatively, PBLUP accounting for inbreeding accumulated since generation 0 while
28 constructing **A** was also used for prediction of breeding values and was named PBLUP_F.

1 The following linear mixed model was used.

$$2 \quad \mathbf{y} = \mathbf{Xb} + \mathbf{Za} + \mathbf{e} \quad (1)$$

3 where \mathbf{y} is a vector of observations; \mathbf{b} is a vector of fixed effects (sex); \mathbf{a} is a vector of direct additive
4 genetic effects of individual animals; \mathbf{e} is a vector of residual errors; and \mathbf{X} and \mathbf{Z} , are known incidence
5 matrices. Assumptions in the model were $\mathbf{a} \sim N(0, \mathbf{A}\sigma_a^2)$ and $\mathbf{e} \sim N(0, \mathbf{I}\sigma_e^2)$; where \mathbf{A} is the numerator
6 relationship matrix between animals derived from pedigree, \mathbf{I} is an identity matrix, and σ_a^2 and σ_e^2 are
7 additive genetic and residual variances, respectively.

8 GBLUP analysis was used to estimate genomic breeding values. A genomic relationship matrix (\mathbf{G}) was
9 constructed by method of Yang et al. (2010) using PLINK 1.9 (Chang et al., 2015). The genomic analysis
10 was done using MTG2 (Lee and van der Werf, 2016) to predict GEBVs for animals in generation 10. The
11 model used for analysis was

$$12 \quad \mathbf{y} = \mathbf{Xb} + \mathbf{Zg} + \mathbf{e} \quad (2)$$

13 where \mathbf{g} is a vector of additive genetic effects of the individual animal, with $\mathbf{g} \sim N(0, \mathbf{G}\sigma_g^2)$ and other
14 terms are defined as above.

15 SSGBLUP analysis combined pedigree and genomic information in a single step extending the animal
16 model to include marker genotypes. The model for evaluation was:

$$17 \quad \mathbf{y} = \mathbf{Xb} + \mathbf{Zu} + \mathbf{e} \quad (3)$$

18 where, \mathbf{u} is the vector of direct additive genetic effects of individual animals, assumed to be normally
19 distributed with $\mathbf{u} \sim N(\mathbf{0}, \mathbf{H}\sigma_u^2)$. The \mathbf{H} matrix includes both non-genotyped and genotyped individuals
20 (Christensen and Lund, 2010; Aguilar et al. 2010). Other terms are already defined as above. This method
21 did not include inbreeding in the set up of \mathbf{A}^{-1} , however, it considered inbreeding in \mathbf{A}_{22}^{-1} . The \mathbf{H}^{-1}
22 derivation was as under

$$\mathbf{H}^{-1} = \begin{bmatrix} \mathbf{A}^{11} & \mathbf{A}^{12} \\ \mathbf{A}^{21} & \mathbf{A}^{22} + \mathbf{G}^{-1} - \mathbf{A}_{22}^{-1} \end{bmatrix}$$

23 where \mathbf{A}_{22} is a pedigree-based numerator relationship matrix for genotyped animals (Aguilar et al., 2010).

1 The tuning of the \mathbf{H}^{-1} matrix component relating to genotyped individuals was according to $\mathbf{A}^{22} + (\mathbf{G}^{*-1} -$
2 $\mathbf{A}_{22}^{-1})$ with \mathbf{G}^* being $\alpha\mathbf{G} + (1-\alpha)\mathbf{A}_{22}$ and three values of α were compared; 0.95, 0.80 and 0.50. In addition
3 to this, the \mathbf{H} matrix was modified as suggested by Vitezica et al. (2011) that involved adding a small
4 constant to all elements of the \mathbf{G} . The value was derived as the difference between the means for \mathbf{A}_{22} and
5 \mathbf{G} . Another single step method (SSGBLUP_F) was employed where \mathbf{A}^{-1} was correctly constructed taking
6 inbreeding into account. The family of BLUPF90 programs (Misztal, 2008; Aguillar et al. 2014) were
7 used to analyse the data.

8 Accuracy of prediction of breeding values was obtained as the Pearson correlation between TBV and
9 GEBV of all individuals in generation 10. Bias was estimated as the deviation of regression coefficient of
10 TBV on GEBV from unity.

11 **Results**

12 *Validation of the simulated data*

13 We showed that in the simulation the observed value for the effective number of chromosome segments
14 (M_e) agreed with the expected M_e , given N_e , as described by Lee et al. (2017). When using 30
15 chromosomes, each with 1 Morgan long, the observed and expected M_e were similar and the effects of
16 the number of historical generations and mutation rate were low (Table S1). This indicates that the
17 simulated populations should have properties that align with the pre-defined population parameter for N_e ,
18 which are relatively robust towards the number of historical generations and the assumed mutation rate in
19 the simulation. Furthermore, we confirmed that observed and expected prediction accuracy agreed well
20 (Table S1), validating the reported values for prediction accuracy in our simulation.

21 *Increasing number of individuals in reference population increases accuracy of GEBV prediction*

22 A larger proportion of genotyped animals led to a larger sample size in the reference population which
23 increased the accuracy of GBLUP as well as for SSGBLUP prediction of GEBV (Table 1). The increase
24 in the accuracy was approximately linear with sample size across the different scenarios where the
25 number of genotyped males increased from 25% to 100% from the last 4 generations.

1 For the RR scenario and a trait with heritability of 0.3 controlled by 990 QTLs, PBLUP gave an accuracy
2 of 0.48 ± 0.01 . The accuracy of G1000 from GBLUP was similar to that of PBLUP. The accuracy
3 increased with a higher sample size *e.g.* 0.61 for G2000 (Table 1). Accuracy with SSGBLUP was much
4 higher; *i.e.* 0.59 ± 0.01 with SS500, increasing to 0.68 ± 0.01 for SS2000. Using genotype information on
5 all animals (G9550) gave the highest prediction accuracy, as expected (0.79 ± 0.01).

6 *Effect of selection on prediction accuracy*

7 In the SR scenario, the accuracy of both EBV and GEBV were lower than in RR. For PBLUP, the
8 accuracy was 0.34 ± 0.01 . A similar value for the prediction accuracy was obtained by GBLUP with
9 G1000 (0.31 ± 0.01). SSGBLUP gave an accuracy of 0.46 ± 0.01 and 0.60 ± 0.01 , with SS500 and SS2000,
10 respectively. For the SA scenario, the prediction accuracy from PBLUP was 0.45 ± 0.01 . A similar
11 prediction accuracy was obtained with G1000 (0.45 ± 0.02) and a much higher accuracy was obtained by
12 SS1000 (0.55 ± 0.02).

13 Variation in α for tuning the contribution from the **G**-matrix to form **H** affected the accuracy of
14 prediction in SSGBLUP for all selection and mating designs. Three values of α were tested; 0.95, 0.80
15 and 0.50. A genetic model with $h^2=0.3$, and 4 QTL models (90, 990, 9990, 60000) were compared. For
16 the RR and SR scenarios, using $\alpha=0.95$ gave the highest accuracy. In the RR scenario, the accuracy
17 decreased by nearly 3.0 to 4.5% and for the SR scenario by 5.0 to 8.0% if α was changed from 0.95 to
18 0.50. However, for the SA scenario, using $\alpha=0.50$ gave the highest accuracy. For the 990-QTL scenario
19 with SA, the SS500 model had the accuracy improved from 0.50 ± 0.03 for $\alpha=0.95$, to 0.51 ± 0.03 for α
20 $=0.80$ and 0.51 ± 0.03 for $\alpha = 0.50$ (Table 3). Improvements from changing α from 0.95 to 0.50 were
21 1.8% for SS1000 and 3.2% for SS2000. For the 90-QTL model, much higher gain in accuracy was
22 observed for the SA scenario by changing α from 0.95 to 0.5 for SS500 (9.8%) and SS1000 (8.7%). For
23 model with a larger number of QTL, the reduction of α did not benefit the accuracy.

24 *Accounting for Inbreeding in A^{-1} increases accuracy of Single Step prediction for populations under* 25 *selection*

26 Accumulation of inbreeding was higher in a selection scenario with assortative mating (Table 5). In the
27 validation population, inbreeding was 2.3% in RR, 3.7% in SR and 9.5% in the SA scenario. For the SA

1 scenario, SSGBLUP (without inbreeding) did not help improve accuracy much over GBLUP. Including
2 inbreeding while constructing \mathbf{A}^{-1} (SSGBLUP_F) resulted in more accurate estimates. For the SS500 we
3 observed a strong improvement in accuracy if inbreeding was accounted for in the genetic evaluation in
4 single step. For the 990QTL model ($h^2=0.3$), accuracy of SS500 increased from 0.50 (SSGBLUP) to 0.57
5 (SSGBLUP_F) (Table 3).

6 *Effect of different heritability estimates and QTL models on prediction accuracy*

7 With $h^2 = 0.1$, the accuracy of prediction was significantly lower across all scenarios for PBLUP, GBLUP
8 and SSGBLUP. For PBLUP, a change of h^2 from 0.3 to 0.5 increased the accuracy by 4.1 to 18.8%, 8.8
9 to 25.0% and 14.0 to 22.7% for RR, SR and SA scenarios, respectively, across the different models.
10 Similarly, a gain in accuracy was observed with increasing h^2 for GBLUP and SSGBLUP across all
11 different genetic models and scenarios. For GBLUP, the magnitude of increase in accuracy with higher
12 heritability was larger, compared to PBLUP and SSGBLUP and more so when a smaller proportion was
13 genotyped. Also for the scenarios with selection (SR and SA) a higher gain in accuracy compared to the
14 RR scenario was observed with higher heritability models.

15 Variation in the number of QTLs did not affect the accuracy of prediction of GEBV for the RR scenario.
16 Similarly, accuracy for EBV obtained by PBLUP was unaffected by variation in QTL model for RR, SR
17 and SA scenarios. Prediction accuracy was affected by QTL number for SR and SA scenarios when using
18 GBLUP and SSGBLUP, where the prediction accuracy improved when the number of QTL in the
19 simulated genetic model increased, with more gains for SA than in the SR scenario (Table 1).

20 *Effect of selective genotyping on the bias of GEBV prediction*

21 In a genetic model with $h^2=0.3$ and 990-QTL, all methods showed no bias with regression coefficients
22 around one for no-selection and random mating (RR) scenario (Table 2). Bias of prediction was a
23 common feature in GBLUP with selective genotyping (G500, G1000 and G2000) in SR and SA design
24 with either under-dispersed or over-dispersed regression coefficients (Table 2). The bias reduced when a
25 larger proportion of males was genotyped, which also led to more genotyped individuals in the reference
26 population. For the SR design, bias in GEBV prediction was evident. G500 and G1000 were most biased
27 with regression coefficients of 0.78 ± 0.06 and 0.85 ± 0.03 , respectively. GBLUP with G2000 was less

1 biased (0.96 ± 0.02). For the SA design, the GBLUP approach gave slightly under-dispersed GEBVs with
2 the regression coefficient equal to 1.11 ± 0.07 and 1.13 ± 0.05 for G500 and G1000, respectively. However,
3 for G2000 and G9550, the estimated regression coefficient was 1.01 ± 0.03 and 0.96 ± 0.01 , respectively,
4 indicating less bias when genomic information is available on more or on all animals.

5 *Effect of selection and mating designs on bias of GEBV prediction*

6 For the SR scenario, the bias of prediction was removed with the SSGBLUP under the 990-QTL model
7 but some bias was observed under the 90-QTL model. However, bias with SSGBLUP was significantly
8 smaller than with GBLUP, especially when only a highly selected proportion was genotyped. For the SA
9 design, the bias in the prediction of breeding value was very high for all methods of prediction, including
10 PBLUP, where a regression coefficient of 0.88 ± 0.02 was obtained. The SSGBLUP methods using $\alpha=0.95$
11 was more biased in the SA scenario. For SS500, regression was 0.68 ± 0.05 , for SS1000 it was 0.72 ± 0.05
12 and even for SS2000 model, the regression was 0.70 ± 0.04 and bias was larger with SSGBLUP than with
13 GBLUP.

14 Reducing the weight given to the **G**-matrix (α from 0.95 to 0.80 and then to 0.50) in SSGBLUP,
15 increased the regression coefficient for the SA scenario, i.e. there was less bias. The improvement in
16 regression coefficient for the 990-QTL ($h^2=0.3$) model by shifting α from 0.95 to 0.50 was 27.9% for
17 SS500, 25% for SS1000 and 24.3% for SS2000 (Table 3). One interesting feature with the SA scenario is
18 that by decreasing the weight of **G** matrix (α from 0.95 to 0.5) in SSGBLUP method also improved the
19 accuracy (Table 3), although the difference was statistically non-significant. The SR scenario showed
20 almost no reduction in accuracy when lowering the value for α .

21 *Including inbreeding in A^{-1} for SSGBLUP reduces bias of prediction of GEBV in population under* 22 *selection*

23 Bias of prediction for the SA scenario was high, even with the SSGBLUP method. However, the bias was
24 practically removed when accounting for inbreeding in the **A** matrix (SSGBLUP_F) (Table 3). For the
25 RR scenario, all estimates were nearly unbiased and inclusion of inbreeding in A^{-1} did not affect the
26 GEBV prediction. For the SR scenario, over-dispersed GEBV (regression coefficients >1) were obtained
27 for SSGBLUP($\alpha=0.50$), however, SSGBLUP_F resulted in regression coefficients close to 1, and at par

1 with PBLUP. For the SA scenario, a tremendous improvement in bias of prediction was observed when
2 inbreeding was accounted for in \mathbf{A} . For the 90QTL scenario, regression coefficient was 0.81 for SS500
3 (SSGBLUP_F) as compared to 0.69 for SS500($\alpha=0.50$), 0.72 for PBLUP and 0.76 for PBLUP including
4 inbreeding (PBLUP_F). This improvement was mainly due to incorporation of inbreeding while
5 constructing \mathbf{A}^{-1} . PBLUP estimations without including inbreeding were similar to PBLUP including
6 inbreeding, although little gain in regression was observed for SA scenario that was in-significant. It is
7 therefore important for the SSGBLUP method to incorporate inbreeding than for the BLUP method,
8 because \mathbf{G}^{-1} will automatically account for inbreeding, and to be consistent with \mathbf{A} , \mathbf{A}^{-1} needs to account
9 for it as well while constructing \mathbf{H}^{-1} .

10 *Effect of heritability and QTL number variation on bias of prediction*

11 With a higher heritability (increasing from 0.3 to 0.5), more accurate predictions with lower bias were
12 obtained across prediction methods and mating designs. However, with a low heritability ($h^2 = 0.1$), the
13 prediction was significantly biased for GBLUP or Single Step methods (Table 2). For the RR scenario,
14 whether using a small or large number of QTLs, there was not much bias observed. However, when
15 selection was involved (SR and SA scenarios), the 90-QTL model resulted in significantly biased
16 predictions for all the methods, including PBLUP and G9550 while models with higher numbers of QTL
17 resulted in lower bias of prediction (Table 2). This happened mainly due to selection of QTL with large
18 effect.

19 **Discussion**

20 In the pedigree-based genetic evaluation of a breeding program, it is assumed that the individuals in the
21 base population are unselected and unrelated having average inbreeding coefficient of zero (Falconer,
22 1996). Henderson (1975) showed that under the infinitesimal genetic model all subsequent selection is
23 conditional upon the unselected base population and is accounted for in BLUP prediction, provided all
24 data is included in the evaluation model. When genomic prediction is based on the genotyped animals
25 alone, this condition is not met, and this gives rise to biased predictions as was shown in the GBLUP
26 scenarios in this study. The bias was manifested in an over-dispersion of the GEBVs. In other words,
27 GEBVs are over-predicted for the top animals at the moment of selection. Bias was smaller when

1 selection for genotyping was less intense and bias was removed when all data that was used in selection
2 decisions was included, as was the case in the SSGBLUP method.

3 Results in our study showed that all prediction methods had lower accuracy in the scenarios with
4 selection (SR and SA). Even the PBLUP method showed a significant (20-30%) decrease in accuracy for
5 all selection scenarios. The lower accuracy in the SR and SA scenarios compared to RR is likely due to
6 the Bulmer effect, *i.e.* the variation between families is reduced due to selection, leading to a reduction of
7 genetic variance and a lower correlation between EBV and TBV (Bijma, 2012). This effect will be
8 relatively large in our study, where the accuracy was validated in the animals from the last generation that
9 had no phenotypic records themselves and the EBV was largely based on information through either
10 pedigree or genomic relationships. The decrease in accuracy was lower for the SSGBLUP compared with
11 the PBLUP scenarios, as the GEBV is more based on information within families (Clark et al. 2013).
12 Assortative mating increases the variance among offspring and shows therefore higher correlations
13 between EBV and TBV. GBLUP prediction accuracy is likely less affected by the Bulmer effect.
14 However, the accuracy in GBLUP was negatively affected by the effect of bias due to only using
15 genotype information on selected animals. With a larger reference population the effect of selective
16 genotyping is smaller, but also leads to more information being used from genotyped individuals leading
17 to more of the within family information being captured and less selection bias.

18 The lower accuracies found in present study for GBLUP compared to SSGBLUP were also reported by
19 Vitezica et al. (2011) and can be attributed to SSGBLUP also using information on un-genotyped
20 individuals that are linked through the pedigree. Comparing G2000 and SS2000 results, we found that the
21 gain in the accuracy was small, which similarly was observed across different scenarios. This was
22 probably because most of the information was already captured by **G** consisting of the individuals from
23 the last 4 generations. Although more deep pedigree information is used in the SS2000 model, the
24 ancestral coefficient of relationship used is numerically very small, giving not much gain in accuracy.
25 However, omitting relationships to the unselected base population from the analysis can still have a
26 significant effect on the ability to correct for selection bias, as was demonstrated by Van der Werf and De
27 Boer (1990).

1 In the present study, data was generated from a population consisting of 50 sires mated to 500 dams, each
2 dam producing 2 progeny. Thus individuals have information on 18 half-sibs and one full-sib, resulting in
3 an estimated N_e -value for the reference population of approximately 358 (Lee et al. 2017). With an N_e -
4 value of 358, the expected value for prediction accuracy derived from theory (Lee et al. (2017))
5 approximately agreed with the observed accuracies in this simulation study (Table S2). For example, the
6 expected accuracies for the scenarios with $h^2=0.3$, and with 50000 SNP according to Lee et al. (2017)
7 were 0.385 (G500), 0.508 (G1000), 0.641 (G2000) and 0.877 for G9550. The observed accuracy (under
8 the no-selection scenario) as obtained with GBLUP for the 990-QTL model were 0.38 ± 0.01 (G500),
9 0.49 ± 0.01 (G1000), 0.61 ± 0.01 (G2000) and 0.79 ± 0.01 (G9550). It is noted that the population structure
10 is more complicated than the full- and half-sib relationships within one generation, which may explain
11 the small difference between the observed and expected prediction accuracy.

12 Accuracy for GEBV obtained by GBLUP was affected by the assumed QTL model. This was probably
13 due to the fact that the gamma distribution employed in the simulation resulted in a few QTL with large
14 effect. In the validation population (Generation 110) the 90-QTL model ($h^2=0.3$) had on average $5.12 \pm$
15 0.35 QTL explaining more than 5% variance individually for the RR scenario. This number was $2.68 \pm$
16 0.28 for SR and 1.6 ± 0.26 for the SA scenario (Table S3). For the 990-QTL model, QTL explaining
17 more than 5% variance were 0.36 ± 0.11 for RR; 0.2 ± 0.08 for SR and 0.16 ± 0.07 for the SA scenario.
18 Two important things are inferred from above data. First, the number of large QTL is larger in the 90-
19 QTL model as compared to the models with ≥ 990 -QTL, mainly due to sharing of TBV over a large
20 number of QTL in the latter models. Second, there is a loss of QTL with large effects in selection
21 scenarios as compared to RR, thus reducing accuracy in selection scenarios. Selection scenarios reduced
22 the genetic as well as phenotypic variance in validation population as compared to first generation (Table
23 6), however, it was seen that the reduction in variance was significantly higher for the 90-QTL model
24 compared to higher QTL-models. For $h^2=0.3$, the loss of genetic variance was 40% for SR and 60% for
25 the SA scenario in 90-QTL model, whereas it was 24% and 36.1% for the 990-QTL model and the
26 reduction was 16.7% and 28.6% for the 60000-QTL model. Loss of a few large loci due to selection
27 might therefore affect the accuracy as well as bias in 90-QTL model significantly.

1 Accuracy was also affected by allele frequency changes, which were observed from generation 101 to
2 generation 110 mostly in selection scenarios (SR and SA). Allele frequency changes were very high for
3 the 90-QTL model. For the 990-QTL model, allele frequency changed but to a lesser extent. Allele
4 frequencies of the five largest QTL (Table S4) revealed that for the 90-QTL model the change of allele
5 frequency in the selection scenario was very high for QTL with large effect. Thus, in SSGBLUP, it
6 negatively affected accuracy in smaller reference, as **A** does not accommodate allele frequency changes.

7 Genomic selection exploits more of the Mendelian sampling variance because it used realized rather than
8 expected relationships (Goddard and Hayes, 2007; Clark et al, 2013). Figure 1 shows that genomic
9 relationships have a normal distribution and there are many negative values. The expected relationships
10 based on pedigree have a skewed distribution with only positive values. Table 4 shows that nearly 50% of
11 the elements in **A** with zero or near zero values are actually negative relationships in the **G** matrix. This
12 helps to explain why pedigree based prediction is less accurate than genomic based predictions. The
13 G9550 have genomic information for all the pedigree that also include the base allelic frequency. This
14 may be the reason that the G9550 scenario gave the most accurate and unbiased prediction of GEBV.
15 Sorensen and Kennedy (1984) and Kennedy et al., (1988), emphasized the fact that the covariances of
16 TBVs for selected individuals are not described well enough by **A** (or **G** for instance with GBLUP),
17 unless all records used in selection are accounted for, as happens for pedigree based estimations. As most
18 of the genomic selection programmes use a **G**-matrix that is obtained from individuals from recent
19 generations, the GEBV estimates are usually biased mostly due to omitting information on selection.

20 Using only selected animals for a reference population can cause biased estimates in GBLUP. The Single
21 Step approach using a **H**-matrix that combines information on **A** and **G** allows to obtain more accurate
22 and less biased estimates of GEBV. The method generally removed the bias seen in GBLUP, although
23 this was not always the case for the SA scenario. Incompatibility of **A** and **G** due to different bases is a
24 thing of concern. In the selection scenario, where highly selected parents or relatives are chosen for
25 constructing reference, the base population frequencies are usually not traceable. We observed in
26 SSGBLUP that by keeping $\alpha=0.95$ and by using the methods of Vitezica et al. (2011), the bias still
27 existed for the SA scenario. Vitezica et al. (2011) proposed a modification for tuning of **H** matrix for
28 populations under selection that involved fitting a constant to all elements of the **G** matrix that they

1 derived by equating the sum of the elements of the \mathbf{G} to the sum of the elements of the \mathbf{A} . Hsu et al.
2 (2017) showed that under selection, if genotypes (SNPs) include QTLs, accuracy and bias of genomic
3 prediction is compromised for Single Step unless the mean of unselected individuals is fitted in the model
4 as a fixed effect. However if the observed SNPs are only markers, the accuracy of prediction may not be
5 improved by the modification proposed by Vitezica et al. (2011). In our data none of the markers (SNPs)
6 were also QTL, although they were generally in LD with QTL.

7 For the large group of non-genotyped animals, breeding values are, a priori, conditioned on genetic
8 values of genotyped animals (Legarra et al. 2009) that are actually based on current genotypic
9 frequencies of recent generations of selected animals, where significant changes in allelic frequency took
10 place due to selection and assortative mating design. It seems likely that the bias with the SSGBLUP
11 methods is caused by an inconsistent scaling of the \mathbf{A} and \mathbf{G} , due to changing allele frequencies with
12 selection. The accumulation of inbreeding in the populations under selection has an effect on the
13 relationship structure in the population. Not all large scale pedigree based genetic evaluation programmes
14 account for inbreeding as it has non-significant influence on the estimations (Meharabani-Yeganeh et al.
15 2000). However, not accounting for inbreeding when deriving relationship will have an effect on the
16 scaling of \mathbf{A} versus \mathbf{G} , genomic relationships automatically account for inbreeding, and this lack of
17 correct scaling can lead to bias of prediction. An inappropriate scaling of \mathbf{G} versus \mathbf{A} is also evident from
18 the decreased bias that was observed when the α -value was decreased from 0.95 to 0.5 (Table 3), i.e.
19 when the genomic relationships are given less weight. However, the bias was removed when inbreeding
20 was incorporated while constructing \mathbf{A}^{-1} . Poor performance of SSGBLUP compared to SSGBLUP_F was
21 therefore due to the presence of highly inbred individuals in data but ignoring inbreeding in the prediction
22 model (Meharabani-Yeganeh et al. 2000, Garcia-Baccino et al. 2017). SSGBLUP is more sensitive to this
23 than PBLUP, which is likely due to the need to combine a relationship matrix \mathbf{G} that accounts for
24 inbreeding implicitly with a pedigree derived matrix that may not account for it. Similarly, BLUP is
25 relatively robust against deviations from the infinitesimal model (Maki-Tanila and Kennedy 1986), but
26 when \mathbf{G} and \mathbf{A} need to be combined, the effect of allele frequency changes seem to be larger. We
27 observed more bias for models with fewer QTL where allele frequency changes are more pronounced.

28 **Conclusion**

1 This study showed that genomic selection using only highly selected genotyped individuals in the
2 reference for genomic prediction results in considerable bias. The Single Step approach resulted in more
3 accurate and less biased estimates of breeding value because it also takes into account the information
4 from non-selected and non-genotyped individuals. However, with selection and assortative mating, some
5 bias was also observed with the Single Step method, likely due to inappropriate merging of the **A** and **G**-
6 matrices due to allele frequency changes of large QTL as a result of selection and also due to ignoring
7 inbreeding in building \mathbf{A}^{-1} . We conclude therefore that the Single Step approach can easily cause bias as
8 it is quite sensitive to inappropriate scaling of **A** and **G**-matrices, especially with selection and selective
9 genotyping, and with considerable rates of inbreeding, but bias can be minimized when scaling is
10 appropriate.

11 **Acknowledgements**

12 GRG duly acknowledge the support provided by Indian Council of Agricultural Research (ICAR) and
13 University of New England Australia, and also financial support provided by Endeavour Research
14 Fellowship (Australia). Dr. Rohan L. Fernando (Iowa State University, USA) and Dr. Andres Legarra
15 (INRA, France) are gratefully acknowledged for the useful discussions. SHL is an Australian Research
16 Council Future Fellow (FT160100229).

17 **Data availability statement**

18 The authors affirm that all data necessary for confirming the conclusions of the article are present within
19 the article, figure, and tables.

20 **References**

- 21 1. Aguilar, I., I. Misztal, S. Tsuruta, A. Legarra & H. Wang, 2014. PREGSF90 – POSTGSF90:
22 Computational Tools for the Implementation of Single-step Genomic Selection and Genome-wide
23 Association with Ungenotyped Individuals in BLUPF90 Programs. Proceedings, 10th World Congress of
24 Genetics Applied to Livestock Production
- 25 2. Bijma, P., 2012. Accuracies of estimated breeding values from ordinary genetic evaluations do not reflect
26 the correlation between true and estimated breeding values in selected populations. *J Anim Breed Genet.*
27 5:345-358.

- 1 17. Lee, S. H., Weerasinghe, S. P., Wray, N. R., Goddard, M. E. and van der Werf, J. H. 2017. Using
2 information of relatives in genomic prediction to apply effective stratified medicine. *Scientific Reports* 7,
3 article number: 42091. doi:10.1038/srep42091 (2017).
- 4 18. Lee, S.H., Clark, S., van der Werf, J.H.J. 2017. Estimation of genomic prediction accuracy from reference
5 populations with varying degrees of relationship. *PLOS ONE* 12, e0189775
- 6 19. Legarra, A., I. Aguilar, & I. Misztal, 2009. A relationship matrix including full pedigree and genomic
7 information. *J. Dairy Sci.* 92: 4656–4663.
- 8 20. Maki-Tanila, A. & B.W. Kennedy, 1986. Mixed-model methodology under genetic models with a small
9 number of additive and non-additive loci. *Proc. 3rd World Congr. Genet. Appl. Livest. Prod.*, vol. 12,
10 p. 443
- 11 21. Mehrabani-Yeganeh, H., Gibson, J. P., Schaeffer, L. R., 2000. Including coefficients of inbreeding in
12 BLUP evaluation and its effect on response to selection. *J. Anim. Breed. Genet.* 117, 145–151
- 13 22. Meuwissen T.H.E., B. J. Hayes & M. E. Goddard, 2001. Prediction of total genetic value using genome-
14 wide dense marker maps. *Genetics*, 157: 1819–1829
- 15 23. Misztal, I. 2008. BLUPF90 - a flexible mixed model program in Fortran 90. *Animal and Dairy Science*,
16 University of Georgia, August 2008.
- 17 24. Patry, C. & V. Ducrocq, 2011. Evidence of biases in genetic evaluations due to genomic preselection in
18 dairy cattle. *J. Dairy Sci.* 94: 1011-1020.
- 19 25. Sargolzaei, M. & F. S. Schenkel, 2009. QMSim: a large-scale genome simulator for livestock.
20 *Bioinformatics* 25 (5): 680-681.
- 21 26. Sorensen, D. A. & B. W. Kennedy, 1984. Estimation of response to selection using least squares and mixed
22 model methodology. *J. Anim. Sci.* 58: 1097–1103.
- 23 27. Van der Werf, J.H.J., and I.J.M. de Boer. 1990. Estimation of additive genetic variance when base
24 populations are selected. *J. Animal Sci.* 68: 3124-3132.
- 25 28. Van Raden, P. M. (2008). Efficient methods to compute genomic predictions. *Journal of Dairy Science* 91,
26 4414–4423.
- 27 29. Van Raden, P. M., C. P. Van Tassell, G. R. Wiggans, T. S. Sonstegard, R. D. Schnabel, J. F. Taylor & F. S.
28 Schenkel, 2009a. Invited review: Reliability of genomic predictions for North American Holstein bulls. *J.*
29 *Dairy Sci.* 92: 16–24.
- 30 30. Van Raden, P. M., M. E. Tooker, & J. B. Cole, 2009b. Can you believe those genomic evaluations for
31 young bulls? *J. Dairy Sci.* 92(E-Suppl. 1): 314 (Abstr.)

- 1 31. Vitezica, Z. G., I. Aguilar, I. Misztal, & A. Legarra, 2011. Bias in genomic predictions for populations
2 under selection. *Genet Res (Camb)* 93: 357–366.
- 3 32. Yang, J., B. Benyamin, B. P. McEvoy, S. Gordon, A. K. Henders, D. R. Nyholt, P. A. Madden, A. C.
4 Health, N. G. Martin, G. W. Montgomery, M. E. Goddard, P. M. Visscher, 2010. Common SNPs explain a
5 large proportion of the heritability for human height. *Nature Genet* 42:565-571.
- 6

Table 1. Estimated accuracy of (G)EBV prediction using different methods and scenarios

heritability	QTLs	Selection and Mating Design	PBLUP	G500	G1000	G2000	SS500	SS1000	SS2000	G9550	
0.10	90	RR	0.35	0.21	0.33	0.41	0.42	0.45	0.49	0.62	
		SR	0.23	0.07	0.12	0.35	0.30	0.34	0.41	0.50	
		SA	0.29	0.12	0.19	0.40	0.20	0.19	0.23	0.51	
	990	RR	0.37	0.22	0.34	0.44	0.45	0.48	0.53	0.64	
		SR	0.26	0.06	0.11	0.36	0.31	0.35	0.45	0.58	
		SA	0.32	0.15	0.23	0.45	0.23	0.26	0.27	0.62	
	60000	RR	0.35	0.22	0.33	0.44	0.45	0.48	0.52	0.64	
		SR	0.26	0.08	0.11	0.37	0.36	0.40	0.48	0.61	
		SA	0.31	0.14	0.23	0.45	0.28	0.31	0.35	0.64	
	Standard Error			0.01-0.02	0.02	0.01-0.02	0.01-0.03	0.01-0.05	0.01-0.04	0.01-0.05	0.01
	0.30	90	RR	0.49	0.38	0.50	0.62	0.59	0.63	0.68	0.78
			SR	0.33	0.22	0.28	0.49	0.44	0.48	0.56	0.63
			SA	0.43	0.32	0.40	0.56	0.41	0.46	0.57	0.65
990		RR	0.48	0.38	0.49	0.61	0.59	0.63	0.68	0.79	
		SR	0.34	0.24	0.31	0.52	0.46	0.51	0.60	0.73	
		SA	0.45	0.39	0.45	0.62	0.50	0.55	0.62	0.77	
60000		RR	0.49	0.37	0.48	0.60	0.58	0.62	0.67	0.78	
		SR	0.32	0.27	0.32	0.54	0.48	0.52	0.62	0.75	
		SA	0.45	0.40	0.47	0.63	0.51	0.59	0.68	0.79	
Standard Error			0.01	0.01-0.02	0.01-0.02	0.01-0.02	0.01-0.04	0.01-0.04	0.01-0.03	0.00-0.01	
0.50		90	RR	0.51	0.44	0.57	0.69	0.64	0.68	0.74	0.85
			SR	0.37	0.27	0.33	0.53	0.46	0.50	0.59	0.65
			SA	0.49	0.41	0.47	0.62	0.54	0.58	0.65	0.68
	990	RR	0.57	0.47	0.59	0.70	0.67	0.71	0.76	0.86	
		SR	0.37	0.34	0.40	0.61	0.54	0.59	0.68	0.81	
		SA	0.54	0.51	0.56	0.71	0.63	0.66	0.75	0.83	
	60000	RR	0.54	0.46	0.58	0.69	0.66	0.70	0.75	0.85	
		SR	0.40	0.35	0.41	0.62	0.55	0.59	0.69	0.83	
		SA	0.54	0.53	0.60	0.73	0.63	0.70	0.77	0.86	
	Standard Error			0.01-0.02	0.01-0.02	0.01-0.02	0.01-0.02	0.01-0.02	0.00-0.03	0.00-0.02	0.00-0.01

RR: No selection and random mating design; SR: Selection on the basis of EBV and random mating design; SA: Selection on the basis of EBV and assortative mating design

PBLUP: Pedigree (10550 with 9550 pedigree and 1000 validation population) based best linear unbiased prediction; G500: GBLUP with 25% close male relatives (125) from 6th to 9th generations in reference; G1000: GBLUP with 50% close male relatives (250) from 6th to 9th generations in reference; G2000: GBLUP with 100% close male relatives (500) from 6th to 9th generations in reference; SS500: Single-Step GBLUP (with $\alpha=0.95$) with pedigree information from 9550 relatives from 9 preceding generations and 25% close male genotyped relatives (125) from 6th to 9th generations in G matrix; SS1000: Single-Step GBLUP (with $\alpha=0.95$) with pedigree information from 9550 relatives from 9 preceding generations and 50% close male genotyped relatives (250) from 6th to 9th generations in G matrix; SS2000: Single-Step GBLUP (with $\alpha=0.95$) with pedigree information from 9550 relatives from 9 preceding generations and 100% close male genotyped relatives (125) from 6th to 9th generations in G matrix; G9550: GBLUP with 9550 relatives from complete pedigree genotyped and used in reference

Table 2. Bias of (G)EBV prediction using different methods and scenarios

heritability	QTLs	Selection and Mating Design	PBLUP	G500	G1000	G2000	SS500	SS1000	SS2000	G9550	
0.10	90	RR	0.99	0.65	1.09	0.99	0.93	0.93	0.93	0.99	
		SR	0.88	1.49	1.45	0.93	0.81	0.82	0.77	0.82	
		SA	0.76	0.18	0.51	0.83	0.26	0.25	0.24	0.77	
	990	RR	1.08	0.77	1.15	1.04	1.00	1.01	1.00	1.02	
		SR	1.08	-0.1	0.37	0.90	0.88	0.92	0.89	0.99	
		SA	0.88	0.82	1.00	0.97	0.32	0.33	0.32	0.94	
	60000	RR	1.02	2.31	1.40	1.09	1.00	1.01	0.99	1.01	
		SR	1.14	0.60	0.46	0.87	1.05	1.06	0.96	1.02	
		SA	0.92	0.38	0.82	0.97	0.46	0.46	0.41	0.98	
	Standard Error			0.03-0.05	0.10-1.61	0.07-2.20	0.03-0.13	0.03-0.07	0.02-0.07	0.02-0.06	0.01-0.02
	0.30	90	RR	1.04	0.97	1.03	1.02	1.00	0.99	1.00	1.00
			SR	0.87	0.67	0.73	0.84	0.85	0.86	0.80	0.80
			SA	0.72	0.81	0.96	0.90	0.49	0.53	0.56	0.78
990		RR	0.95	1.16	1.01	1.01	1.01	1.00	1.00	1.00	
		SR	1.00	0.78	0.85	0.96	1.00	1.02	0.96	0.98	
		SA	0.88	1.11	1.13	1.01	0.68	0.72	0.70	0.96	
60000		RR	0.99	1.02	1.00	1.00	0.97	0.97	0.97	1.01	
		SR	0.93	0.84	0.87	0.96	1.02	1.04	0.97	1.00	
		SA	0.91	1.26	1.17	1.04	0.73	0.79	0.77	0.99	
Standard Error			0.01-0.03	0.04-0.12	0.03-0.06	0.01-0.03	0.01-0.06	0.01-0.05	0.01-0.03	0.01-0.02	
0.50		90	RR	0.93	0.97	0.98	0.98	0.95	0.96	0.97	0.99
			SR	0.75	0.67	0.74	0.84	0.76	0.79	0.75	0.75
			SA	0.78	0.85	0.94	0.87	0.63	0.67	0.67	0.75
	990	RR	1.02	1.07	1.04	1.04	1.00	1.00	1.00	1.01	
		SR	0.91	0.85	0.92	0.97	1.02	1.06	0.97	0.98	
		SA	0.93	1.09	1.13	1.01	0.86	0.88	0.85	0.96	
	60000	RR	0.99	1.05	1.01	0.99	0.98	0.98	0.98	1.00	
		SR	1.01	0.88	0.96	0.99	1.05	1.10	1.00	1.00	
		SA	0.95	1.15	1.23	1.05	0.83	0.91	0.88	1.00	
	Standard Error			0.01-0.03	0.02-0.05	0.02-0.03	0.01-0.02	0.01-0.05	0.01-0.04	0.01-0.03	0.01-0.02

Footnotes for Table 2 are the same as in Table 1.

Table 3. Estimated accuracy and bias of GEBV prediction in Single Step approach with tuning of H matrix

Criteria	QTLs	Selection and Mating Design	SS500 $\alpha=0.95,$ \$	SS1000 $\alpha=0.95,$ \$	SS2000 $\alpha=0.95,$ \$	SS500 $\alpha=0.80$	SS1000 $\alpha=0.80$	SS2000 $\alpha=0.80$	SS500 $\alpha=0.50$	SS1000 $\alpha=0.50$	SS2000 $\alpha=0.50$	SS500 F	SS1000 F	SS2000 F
Accuracy	90	RR	0.59	0.63	0.68	0.59	0.63	0.68	0.58	0.61	0.66	0.59	0.63	0.68
		SR	0.44	0.48	0.56	0.44	0.48	0.56	0.42	0.45	0.55	0.44	0.48	0.57
		SA	0.41	0.46	0.57	0.43	0.49	0.57	0.45	0.50	0.58	0.51	0.55	0.61
	990	RR	0.59	0.63	0.68	0.59	0.63	0.68	0.57	0.61	0.65	0.59	0.63	0.68
		SR	0.46	0.51	0.60	0.46	0.50	0.60	0.43	0.47	0.57	0.46	0.51	0.60
		SA	0.50	0.55	0.62	0.51	0.55	0.64	0.51	0.56	0.64	0.57	0.61	0.60
	60000	RR	0.58	0.62	0.67	0.58	0.62	0.66	0.57	0.60	0.64	0.58	0.62	0.67
		SR	0.48	0.52	0.62	0.47	0.51	0.61	0.44	0.48	0.59	0.48	0.52	0.62
		SA	0.51	0.59	0.68	0.53	0.59	0.68	0.57	0.60	0.67	0.60	0.63	0.62
	Standard Error			0.01- 0.04	0.01- 0.04	0.01- 0.03	0.01- 0.02	0.01- 0.02	0.01- 0.02	0.01- 0.03	0.01	0.01- 0.02	0.01- 0.02	0.01- 0.02
Bias	90	RR	1.00	0.99	1.00	1.04	1.03	1.04	1.10	1.11	1.13	1.03	1.01	1.02
		SR	0.85	0.86	0.80	0.90	0.92	0.87	0.96	1.00	1.00	0.90	0.91	0.86
		SA	0.49	0.53	0.56	0.54	0.58	0.60	0.65	0.69	0.70	0.81	0.82	0.81
	990	RR	1.01	1.00	1.00	1.04	1.05	1.04	1.09	1.11	1.23	1.03	1.03	1.02
		SR	1.00	1.02	0.96	1.05	1.08	1.03	1.12	1.16	1.18	1.06	1.08	1.02
		SA	0.68	0.72	0.70	0.74	0.78	0.76	0.87	0.90	0.87	1.03	1.04	1.02
	60000	RR	0.97	0.97	0.97	1.00	1.01	1.01	1.05	1.07	1.09	0.99	1.00	0.99
		SR	1.02	1.04	0.97	1.08	1.11	1.05	1.16	1.20	1.21	1.09	1.10	1.03
		SA	0.73	0.79	0.77	0.79	0.84	0.82	0.95	0.97	0.94	1.08	1.09	1.03
	Standard Error			0.02- 0.06	0.01- 0.05	0.01- 0.04	0.02- 0.04	0.01- 0.04	0.01- 0.03	0.02- 0.04	0.02- 0.04	0.01- 0.04	0.01- 0.03	0.01- 0.04

The simulated data with $h^2 = 0.3$ were used; \$: modification of SSGBLUP as suggested by Vitezica et al. (2011), $\alpha=0.50$: shrinking information from **G** matrix to 50%, $\alpha=0.95$: shrinking information from **G** matrix to 95%, F: SSGBLUP_F accounting for inbreeding in \mathbf{A}^{-1} and \mathbf{A}_{22}^{-1} .

Table 4. The distribution of pair-wise relationships for the individuals from the GRM and NRM in different selection and mating designs

#Genetic relationship	RR		SR		SA	
	GRM	NRM	GRM	NRM	GRM	NRM
-0.25 to 0.00	302,20,919	141,09,673	303,08,050	144,79,933	305,74,102	147,68,775
0.001 to 0.25	253,25,866	414,06,656	252,37,460	410,32,352	249,62,864	403,93,558
0.251 to 0.50	87,462	1,13,318	88,798	1,16,840	96,957	4,59,724
0.501 to 0.75	11,711	16,313	11,656	16,845	12,006	23,520
0.751 to 0.85	17	15	11	5	44	335
0.851 to 0.95	0	0	0	0	2	58
>0.95 (up to 0.97)	0	0	0	0	0	5

The total number of pair-wise relationships was 55645975

RR: No selection and random mating design; SR: Selection on the basis of EBV and random mating design; SA: Selection on the basis of EBV and assortative mating design

GRM: genomic relationship matrix, NRM: numerator relationship matrix

#All values for NRM in -0.25 to 0.00 category are actually belong to 0.00 and no value is negative in NRM

Table 5. The distribution of inbreeding coefficients in generations for different selection and mating design scenarios

Generation	Individuals	Inbreeding		
		RR	SR	SA
Founders	550	0.000	0.000	0.000
101	1000	0.000	0.000	0.000
102	1000	0.002	0.003	0.003
103	1000	0.006	0.006	0.009
104	1000	0.006	0.011	0.020
105	1000	0.010	0.014	0.031
106	1000	0.011	0.018	0.039
107	1000	0.015	0.023	0.047
108	1000	0.017	0.026	0.060
109	1000	0.020	0.032	0.072
110(validation)	1000	0.023	0.037	0.095
Min		0.000	0.000	0.000
Max		0.258	0.271	0.344

RR: No selection and random mating design; SR: Selection on the basis of EBV and random mating design; SA: Selection on the basis of EBV and assortative mating design

Table 6. Estimates of phenotypic and genetic variance in the first generation and validation population of the 10th generation

Scenario	*QTL model	h ² =0.1				h ² =0.3				h ² =0.5			
		Phenotypic variance		Genetic variance		Phenotypic variance		Genetic variance		Phenotypic variance		Genetic variance	
		G1 (0.03-0.05)	G10 (0.03-0.05)	G1 (0.01)	G10 (0.01)	G1 (0.03-0.06)	G10 (0.04-0.06)	G1 (0.01-0.03)	G10 (0.02-0.04)	G1 (0.04-0.07)	G10 (0.04-0.08))	G1 (0.03-0.06)	G10 (0.02-0.09)
RR	90	1.02	0.99	0.10	0.10	0.99	0.99	0.30	0.30	1.00	0.95	0.50	0.47
	990	1.00	1.00	0.10	0.10	1.01	1.00	0.30	0.30	0.98	0.99	0.49	0.50
	60000	1.00	0.99	0.10	0.10	0.99	0.99	0.30	0.30	1.00	0.98	0.49	0.48
SR	90	1.00	0.98	0.10	0.07	1.00	0.86	0.30	0.18	0.99	0.72	0.49	0.22
	990	1.02	0.98	0.10	0.08	1.00	0.94	0.29	0.22	0.99	0.87	0.50	0.38
	60000	1.01	1.00	0.10	0.09	1.00	0.95	0.30	0.25	1.00	0.90	0.50	0.39
SA	90	1.00	0.95	0.10	0.06	1.03	0.85	0.35	0.14	1.14	0.70	0.64	0.20
	990	1.00	0.98	0.10	0.08	1.05	0.92	0.36	0.23	1.13	0.88	0.62	0.38
	60000	1.01	0.98	0.10	0.08	1.03	0.94	0.35	0.25	1.14	0.93	0.63	0.39

Numbers in the parentheses are estimates for range of Standard Deviation (SD)

RR: No selection and random mating design; SR: Selection on the basis of EBV and random mating design; SA: Selection on the basis of EBV and assortative mating design

*QTL model: number of QTLs explaining the variance of the trait

G1: Generation 1 (First) with 1000 individuals, G10: Validation population in the 10th Generation (N=1000)

SUPPLEMENTARY TABLES:

Table S1: Validation for the simulated test data with theory (Lee et al. 2017)

Test Data X						
Parameters	Scenario1	Scenario2	Scenario3	Scenario4	Scenario5	Scenario6
Ne	100	100	100	500	500	500
Generations	50	100	200	250	500	1000
Observed Me±SE	246.3±1.54	242.2±1.37	230.4±1.31	1016.8±1.85	988.8±2.00	940.3±1.72
Expected Me	246			1157		
Test Data Y						
Parameters	Scenario1	Scenario2	Scenario3	Scenario4	Scenario5	Scenario6
Ne	100	100	100	500	500	500
Generations	100	100	100	500	500	500
Mutation Rate	2.5×10^{-3}	2.5×10^{-5}	2.5×10^{-8}	2.5×10^{-3}	2.5×10^{-5}	2.5×10^{-8}
Observed Me±SE	272.9±1.38	246.2±1.09	248.3±0.82	1171±2.69	1006.5±0.74	1007.2±1.53
Expected Me	246			1157		
Test Data Z						
Observed Me±SE	Expected Me		Observed accuracy	Expected accuracy		
232.36±1.42	246		0.818±0.003	0.818		

In the test data X, variations included effective population size (Ne) of 100 for historical population with generations 50, 100 and 200, similarly for Ne of 500, historical generations of 250, 500 and 1000 were simulated. 100 replicates of test data X were run to obtain Me. In test data Y, the sensitivity to the mutation rate was analysed for historical populations with three different mutation rates namely high (2.5×10^{-3}), medium (2.5×10^{-5}) and low (2.5×10^{-8}). Ten replicates of test data Y were run to obtain Me. In the test data Z, with 100 Ne and 100 historical generations with mutation rate of 2.5×10^{-8} , in the last generation, the population size was increased to 3000 and out of 3000 genotyped individuals, 2000 were included in the reference and 1000 in validation population. Accuracy was obtained as Pearson's correlation between TBV and GEBV. 35 replications were carried out for test data Z. Thirty chromosomes, each with 1 Morgan long, were used.

Table S2: Validation for the RR design simulated data with the theory (Lee et al. 2017)

Data simulated for main study of the manuscript (RR Design)					
heritability	GBLUP	G500	G1000	G2000	G9550
0.10	Expected accuracy with Ne=358	0.234	0.322	0.434	0.725
	Observed accuracy with GBLUP	0.22±0.02	0.34±0.01	0.44±0.01	0.64±0.01
0.30	Expected accuracy with Ne=358	0.385	0.508	0.641	0.877
	Observed accuracy with GBLUP	0.38±0.01	0.49±0.01	0.61±0.01	0.79±0.01
0.50	Expected accuracy with Ne=358	0.474	0.606	0.733	0.920
	Observed accuracy with GBLUP	0.47±0.01	0.59±0.01	0.70±0.01	0.86±0.00

RR: No selection and random mating design; G500: GBLUP with 25% close male relatives (125) from 6th to 9th generations in reference; G1000: GBLUP with 50% close male relatives (250) from 6th to 9th generations in reference; G2000: GBLUP with 100% close male relatives (500) from 6th to 9th generations in reference; G9550: GBLUP with 9550 relatives from complete pedigree genotyped and used in reference. Observed estimates are provided with standard error. Thirty chromosomes, each with 1 Morgan long, were used.

Table S3: Number of Large QTL in validation population (Generation 10) explaining per cent variance

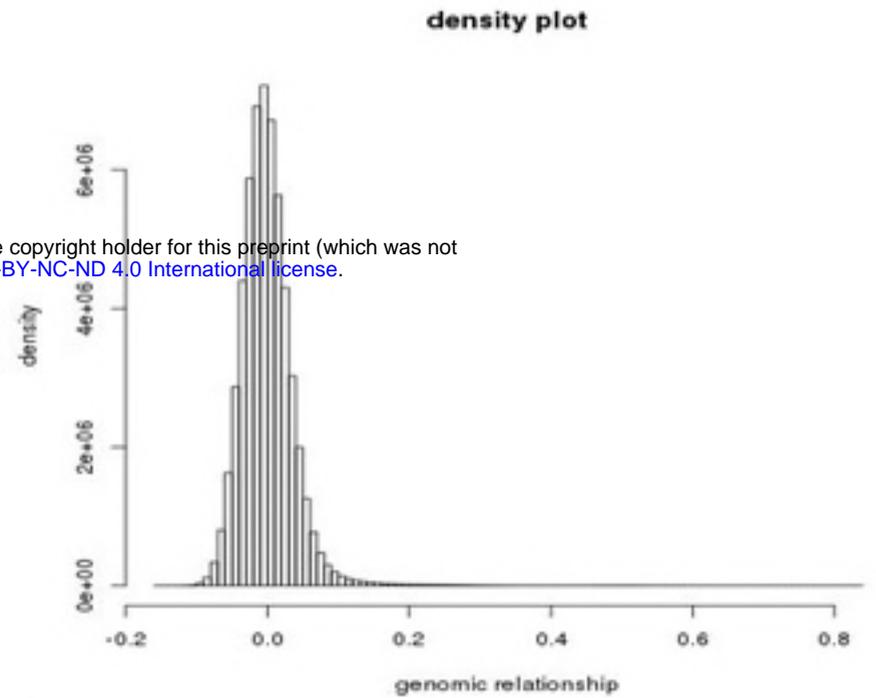
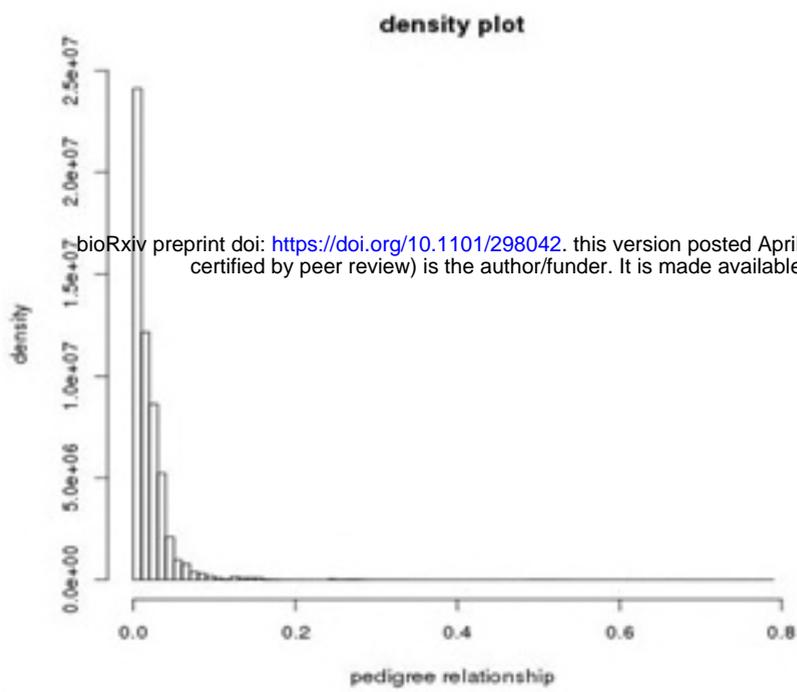
Design ($h^2=0.30$)	QTL-model	1%	5%	10%	25%
RR	90	18.08±0.70	5.12±0.35	1.96±0.16	0.32±0.10
	990	18.6±0.60	0.36±0.11	nil	nil
	60000	Nil	nil	nil	nil
SR	90	15.88±0.72	2.68±0.28	0.68±0.22	0.04±0.04
	990	17.6±0.68	0.20±0.08	0.12±0.08	nil
	60000	Nil	nil	nil	nil
SA	90	13.44±0.62	1.60±0.26	0.12±0.09	nil
	990	13.76±2.75	0.16±0.07	nil	nil
	60000	Nil	nil	nil	nil

Table S4: # Allele frequency changes for genetic model in 10 generations

Design	Genetic model ($h^2=0.30$)	Generation1		Generation10	
		Allele-1	Allele-2	Allele-1	Allele-2
RR	90-QTL model	0.55	0.45	0.61	0.39
		0.51	0.49	0.39	0.61
		0.25	0.75	0.27	0.73
		0.78	0.22	0.82	0.18
		0.53	0.47	0.60	0.40
	990-QTL model	0.63	0.37	0.65	0.35
		0.53	0.47	0.59	0.41
		0.40	0.60	0.44	0.56
		0.44	0.56	0.44	0.56
		0.51	0.49	0.42	0.58
SR	90-QTL model	0.79	0.21	0.17	0.83
		0.05	0.95	0.90	0.10
		0.13	0.87	0.74	0.26
		0.34	0.66	0.15	0.85
		0.47	0.53	0.54	0.46
	990-QTL model	0.41	0.59	0.21	0.79
		0.81	0.19	0.98	0.02
		0.55	0.45	0.67	0.33
		0.89	0.11	0.98	0.02
		0.74	0.26	0.96	0.04
SA	90-QTL model	0.45	0.55	0.99	0.01
		0.80	0.20	1.00	0.00
		0.43	0.57	0.07	0.93
		0.50	0.50	0.24	0.76
		0.80	0.20	0.37	0.63
	990-QTL model	0.48	0.52	0.90	0.10
		0.25	0.75	0.04	0.96
		0.13	0.87	0.62	0.38
		0.71	0.29	0.44	0.56
		0.45	0.55	0.80	0.20

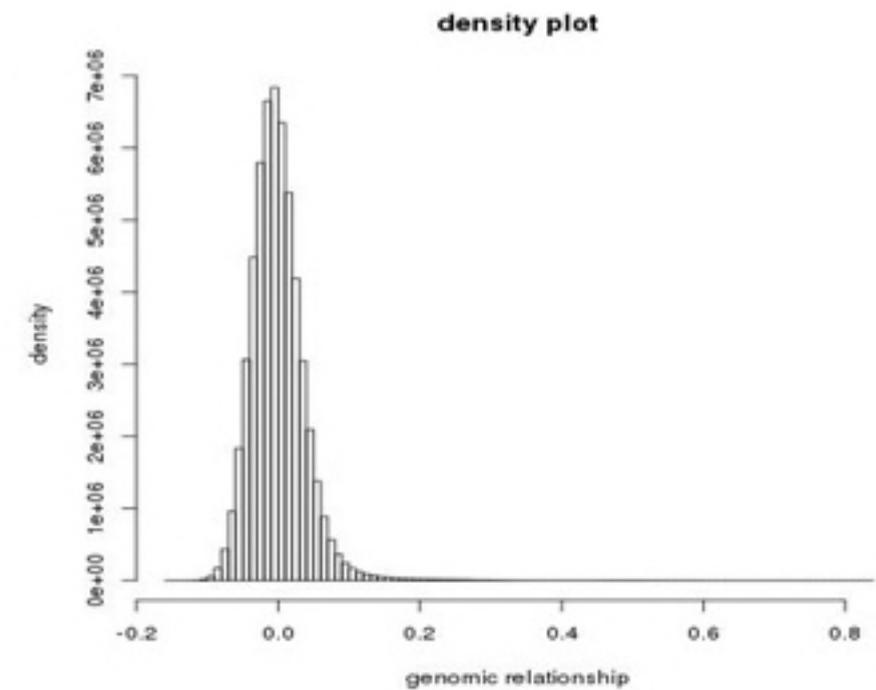
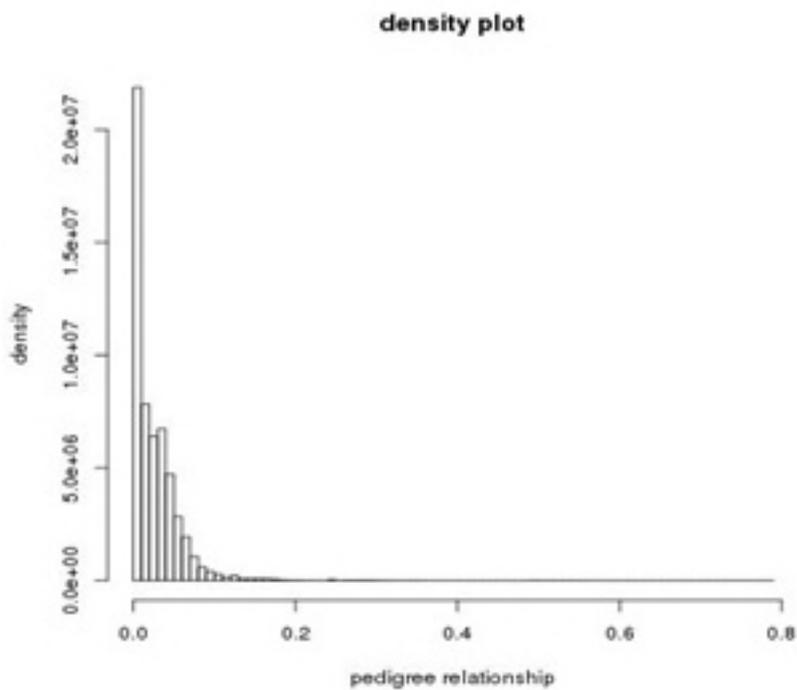
[#]The table describes allele frequency for 5 biggest QTL under several genetic models. As the QTL is bi-allelic, frequency of two alleles for each QTL is shown.

No selection and Random Mating (RR)



bioRxiv preprint doi: <https://doi.org/10.1101/298042>; this version posted April 9, 2018. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. It is made available under a [CC-BY-NC-ND 4.0 International license](https://creativecommons.org/licenses/by-nc-nd/4.0/).

Selection through EBV and random mating (SR)



Selection through EBV and positive assortative mating (SA)

