

Research

Open Access

## Persistence of accuracy of genome-wide breeding values over generations when including a polygenic effect

Trygve R Solberg<sup>\*1</sup>, Anna K Sonesson<sup>2</sup>, John A Woolliams<sup>1,3</sup>,  
Jørgen Ødegard<sup>2,1</sup> and Theo HE Meuwissen<sup>1</sup>

Address: <sup>1</sup>Norwegian University of Life Sciences, Department of Animal and Aquacultural Sciences, PO Box 5003, N-1432 Ås, Norway, <sup>2</sup>NOFIMA Marine, PO Box 5010, N-1432 Ås, Norway and <sup>3</sup>Roslin Institute (Edinburgh), Roslin, Midlothian EH25 9PS, UK

Email: Trygve R Solberg<sup>\*</sup> - trygve.roger.solberg@umb.no; Anna K Sonesson - anna.sonesson@nofima.no;  
John A Woolliams - john.woolliams@bbsrc.ac.uk; Jørgen Ødegard - jorgen.odegard@nofima.no;  
Theo HE Meuwissen - theo.meuwissen@umb.no

<sup>\*</sup> Corresponding author

Published: 29 December 2009

Received: 30 January 2009

Genetics Selection Evolution 2009, **41**:53 doi:10.1186/1297-9686-41-53

Accepted: 29 December 2009

This article is available from: <http://www.gsejournal.org/content/41/1/53>

© 2009 Solberg et al; licensee BioMed Central Ltd.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/2.0>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

### Abstract

**Background:** When estimating marker effects in genomic selection, estimates of marker effects may simply act as a proxy for pedigree, i.e. their effect may partially be attributed to their association with superior parents and not be linked to any causative QTL. Hence, these markers mainly explain polygenic effects rather than QTL effects. However, if a polygenic effect is included in a Bayesian model, it is expected that the estimated effect of these markers will be more persistent over generations without having to re-estimate the marker effects every generation and will result in increased accuracy and reduced bias.

**Methods:** Genomic selection using the Bayesian method, 'BayesB' was evaluated for different marker densities when a polygenic effect is included (GWpEBV) and not included (GWEBV) in the model. Linkage disequilibrium and a mutation drift balance were obtained by simulating a population with a  $N_e$  of 100 over 1,000 generations.

**Results:** Accuracy of selection was slightly higher for the model including a polygenic effect than for the model not including a polygenic effect whatever the marker density. The accuracy decreased in later generations, and this reduction was stronger for lower marker densities. However, no significant difference in accuracy was observed between the two models. The linear regression of TBV on GWEBV and GWpEBV was used as a measure of bias. The regression coefficient was more stable over generations when a polygenic effect was included in the model, and was always between 0.98 and 1.00 for the highest marker density. The regression coefficient decreased more quickly with decreasing marker density.

**Conclusions:** Including a polygenic effect had no impact on the selection accuracy, but showed reduced bias, which is especially important when estimates of genome-wide markers are used to estimate breeding values over more than one generation.

## Background

High-throughput genotyping and availability of dense marker information have made prediction of breeding values based on dense marker genotyping possible resulting in so-called genome-wide breeding values (GWEBV). Several methods have been suggested to estimate marker effects in the prediction of GWEBV e.g. [1-4]. The advantage of selecting parents on GWEBV in genomic selection schemes is the potential to select candidates with high accuracy and low bias directly by using marker genotypes or haplotypes only. Several simulation studies in which markers were calibrated on a training set of phenotypes e.g. [5-8] have demonstrated these advantages. GWEBV may reduce the amount of phenotyping required for breeding schemes and hence constitute an attractive proposition because obtaining phenotypic records routinely for all generations can be expensive, reduce animal welfare, and sometimes impossible for live selection candidates. An example for which all three issues apply is in fish aquaculture where selecting disease resistance is done by challenging sibs of the candidates with the disease to avoid infecting the selection candidates [9]. Genomic selection selects animals directly on the genotype, rather than the phenotype which may be an advantage, especially for traits that cannot be measured or are expensive to record on selection candidates (e.g. slaughter traits and challenge-test data).

Successful implementation of genomic selection relies on some underlying assumptions. One assumption is the existence of population-wide linkage disequilibrium (LD) between markers and quantitative trait loci (QTL). As a result of imperfect LD, markers in LD with the QTL are not likely to explain all existing genetic variation, and the remaining genetic variation will be included in the polygenic variance. For sparse marker maps, linkage disequilibrium between the markers and the QTL will be reduced and only part of the genetic variance will be explained by the markers. Furthermore, estimated marker effects may model family relationships [6], which will result in spurious associations between phenotypes and marker alleles, i.e. there will be non-zero marker effects whilst the marker alleles are not linked to any causative QTL. It is expected that such marker-QTL associations decay at a rate of  $(1-c)^t$ , where  $c$  is the distance between marker and QTL and  $t$  is the number of generations [10]. For spurious associations,  $c = 0.5$ , and for tightly linked markers,  $c < 0.01$ . Thus, spurious associations decay much more rapidly than associations based on real linkage over time. This introduces the important issue of the persistence of GWEBV predictions over generations in the absence of marker effects re-estimation, and few studies have examined this issue, e.g. [11].

One solution that may address both issues is to include a polygenic effect in the model and this has been addressed by others, e.g. [6,7,12], but not evaluated over multiple generations. Our hypothesis is that these spurious associations are better represented by a polygenic effect in the model than by markers that happen to have a higher frequency in some families compared to others. Thus it is expected that including a polygenic effect will capture genetic variation that is not in tight linkage with markers, such that this variation is to a lesser extent captured by markers through spurious associations. The objective of this study was to test this hypothesis by investigating the effect of including a polygenic effect into a Bayesian model for the estimation of marker estimates to predict GWEBV, and their accuracy and bias over multiple generations to study their persistence over time.

## Methods

### Population structure and genome size

Details of the simulation model have been described in an earlier paper [8]. Briefly, a population with an effective population size of  $N_e = 100$  was simulated over 1000 generations of random mating, random selection and with a genome subject to mutation. In generation  $t = 1001$ , the number of animals was increased to 1000 by factorial mating of 50 sires ( $i = 1-50$ ) and 50 dams ( $i = 51-100$ ) from generation  $t = 1000$ . The factorial mating was achieved by mating sire 1 to dams 51-70, sire 2 to dams 52-71, sire 3 to dams 53-72 and so on, and each dam had one offspring per sire. In descending generations ( $t = 1002$  to  $t = 1006$ ), the animals had 1000 offspring produced by random sampling with replacement among the parents selected from the previous generation.

The size and structure of the genome were the same as described in [8]. The genome (10 Morgan) was simulated with 10 chromosomes each 100 cM long. Four density schemes were evaluated, and the density was scaled by the effective population size ( $N_e$ ) used to generate the markers, which was  $N_e = 100$  and a genome size in Morgan (M). Scaled marker densities were 1, 2, 4 and 8  $N_e/M$ , which corresponded here to 100, 200, 400 and 800 markers per Morgan.

Mendelian inheritance and the Haldane mapping function were assumed for all loci. The mutation rate of the markers was assumed to be  $2.5 \times 10^{-3}$  per locus per meiosis. With this mutation rate, 99% of the potential markers were segregating at  $t = 1001$ . Markers with more than two alleles segregating at  $t = 1001$  were converted to bi-allelic SNP markers by ignoring some of the mutations as described in [8]. The minor allele frequencies of the SNP markers showed approximately a uniform distribution with an over-representation of marker alleles with intermediate frequencies, which in practice may reflect the

effect of pre-screening SNP markers and selecting the most informative. The potential number of QTL was kept at 100 per chromosome, distributed evenly over each chromosome. The actual number of segregating QTL at  $t = 1001$  depended on the mutation rate which was assumed to be  $2.5 \times 10^{-5}$  per locus per meiosis. The resulting number of segregating QTL was typically 5 to 6% of the potential number. The additive effect of a mutational allele of the multi-allelic QTL was sampled from the gamma distribution with a shape parameter of 0.4 and scale parameter of 1.66 [13] with an equal probability of a positive or negative effect. No polygenic effect was simulated.

### True breeding value (TBV) and phenotypic values

The true breeding value (TBV) of animal  $i$  from generation 1001-1006 was calculated as:

$$TBV_i = \sum_{j=1}^{N_{QTL}} Q_{ij} q_j$$

where  $q_j$  is a vector of true QTL effects of the QTL alleles at locus  $j$ , and  $Q_{ij}$  is an incidence row vector indicating for animal  $i$  which of the QTL alleles it carried at locus  $j$  (e.g.  $Q_{ij} = [1 \ 1 \ 0 \ \dots]$  for animal  $i$  carrying QTL alleles 1 and 2 at locus  $j$ );  $N_{QTL}$  is the number of QTL loci. For generation 1001, phenotypic values for each animal were simulated as:  $y_i = TBV_i + \varepsilon_i$ , where  $\varepsilon_i \sim N(0, \sigma_e^2)$ . The variance of the TBV effects ( $\sigma_{TBV}^2$ ) varied somewhat from replicate to replicate, but was on average 1.0 (s.e. = 0.118). The environmental variance ( $\sigma_e^2$ ) was set equal to  $\sigma_{TBV}^2$  such that the heritability was 0.5 for every replicate.

### Estimation model with polygenic effect

The 'BayesB' method of Meuwissen *et al.* [3] was used to estimate the effects of the SNP markers for the 1000 animals in generation  $t = 1001$ . The 'BayesB' model is described in more detail in earlier papers [3] and [8], but briefly, the variance of the marker effects ( $\sigma_{gj}^2$ ) was estimated for every marker using a relevant prior distribution which was a mixture of an inverted chi-squared distribution and a discrete probability mass at  $\sigma_{gj}^2 = 0$ . A Metropolis-Hastings algorithm was used to sample  $\sigma_{gj}^2$  from its distribution conditional on  $y^*$ ,  $p(\sigma_{gj}^2 | y^*)$ , where  $y^*$  denotes the data  $y$  corrected for the mean and all other genetic effects except the marker effect ( $g_j$ ) [14]. Given  $\sigma_{gj}^2$ , marker effects,  $g_j$  were sampled from a Normal distribution as prior and using Gibbs sampling [15]. The 'BayesB' model was extended to include a polygenic effect ( $a$ ):

$$y = 1\mu + a + \sum_{j=1}^{N_{loc}} X_j g_j + e$$

where  $y$  is the vector of phenotypes,  $\mu$  is the overall mean,

$a$  is the vector of polygenic effects,  $\sum_{j=1}^{N_{loc}}$  is the summation

over all marker loci from 1 to  $N_{loc}$ , where  $N_{loc}$  is varying from 1010 marker loci for the lowest marker density ( $1Ne/M$ ) to 8080 marker loci for the highest marker density ( $8Ne/M$ ).  $X_j$  is a design matrix for the  $j$ 'th marker,  $g_j$  is the vector of the  $j$ 'th marker effect and  $e$  is the residual term. Dimension of the  $y$ ,  $a$ , and  $e$  vectors are  $1000 \times 1$ , the  $X_j$  matrix varies from  $1010 \times 2$  for the lowest marker density up to  $8080 \times 2$  for the highest marker density. The variance of  $a$  was  $Var(a) = A\sigma_a^2$ , where  $A$  ( $1000 \times 1000$ ) is the additive relationship matrix, calculated based on five generations of pedigree from generation  $t = 996$  to  $t = 1000$  using the algorithm of [16]. Polygenic effects were sampled in the MCMC chain using Gibbs sampling and assuming a prior  $N(0, \sigma_a^2)$  following [15], and  $\sigma_a^2$  was estimated using a scaled inverted chi-squared prior distribution with -2 degrees of freedom, which implies a non-informative flat prior distribution [15].

The variance of the marker effect was  $\sigma_{gj}^2$ , which was estimated for every marker using a mixture distribution as the prior;

$$p(\sigma_{gj}^2) = \begin{cases} 0 & \text{with probability } (1-p) \\ \chi^{-2}(\nu, S) & \text{with probability } p \end{cases}$$

The probability  $p$  depends on the density of the markers, and varies with different marker densities, because with more markers, it becomes less likely for marker  $j$  to be required to capture the predictive LD between QTL and markers, i.e.  $p = 53/(N_{loc})$  where 53 is the expected number of QTL and  $N_{loc}$  is the number of marker loci. Sampling from the posterior distribution of  $\sigma_{gj}^2$  was by a Metropolis-Hastings algorithm that sampled  $\sigma_{gj}^2$  from  $p(\sigma_{gj}^2 | y^*)$ , where the prior distribution,  $p(\sigma_{gj}^2)$ , was used as the distribution to suggest updates for the Metropolis Hastings chain [13], and  $y^*$  denotes the data  $y$  corrected for the mean and all other genetic effects except the marker effect ( $g_j$ ). The Metropolis Hastings chain was run for 10.000 cycles using a burn-in period of 1000 cycles. Given  $\sigma_{gj}^2$ , marker effects,  $g_j$  were sampled from  $p(g_j | \sigma_{gj}^2)$  using Gibbs sampling [15].

### Prediction of genome-wide breeding values

Prediction of the GWEBV for the method 'BayesB' without polygenic effect was calculated from:

$$\text{GWEBV}_i = \sum_{j=1}^{N_{loc}} X_{ij} \hat{g}_j$$

where  $X_{ij}$  denotes the marker genotype of animal  $i$  at locus  $j$  in generation  $t = 1002$  to  $t = 1006$ , and  $\hat{g}_j$  is the estimate of the marker effects, which was estimated on animals in generation  $t = 1001$ .

Prediction of the breeding values including the polygenic effect (GWpEBV) for the method 'BayesB' was calculated from:

$$\text{GWpEBV}_i = \sum_{j=1}^{N_{loc}} X_{ij} \hat{g}_j + \hat{a}_{i(t)}.$$

Since no data was recorded after generation  $t = 1001$ , the polygenic effect  $\hat{a}_{i(t)}$  of animal  $i$  in generation  $t$  was calculated as  $\hat{a}_{i(t)} = 0.5 \hat{a}_{s(t-1)} + 0.5 \hat{a}_{d(t-1)}$  where the subscripts  $s$  and  $d$  represent the sire and dam of animal  $i$ , respectively. This formula is valid here because the parents of the next generation were randomly selected and there was no phenotypic data entering the evaluations in later generations, i.e. after generation  $t = 1001$ . For each replicate, the mean and median of the Gibbs samples for the polygenic variance were calculated from the final 5000

values of the chain. These values were then averaged over 10 replicates.

As a measure of bias we calculated the linear regression coefficient of true breeding values on GWEBV and GWpEBV within each of the five generations from  $t = 1002$  to  $t = 1006$ . The correlation coefficients were calculated between the true breeding value and the GWEBV and GWpEBV for all five generations which reflects the accuracy of predicting the genome-wide breeding values. The result is based on the average of 20 replicates for each marker density.

## Results

### Accuracy of selection

Tables 1 to 4 show the accuracy of selection and bias for the four marker densities when the polygenic effect (GWEBV) is ignored and when it is included (GWpEBV). For the highest marker density (8Ne/M), the selection accuracy decreased from 0.875 in generation  $t = 1002$  to 0.842 in generation  $t = 1006$  for GWEBV (Table 1). Selection accuracy was higher for GWpEBV than for GWEBV for all generations and this was particularly significant for three out of five generations. The difference in accuracy between the two models varied from 0.008 in generation  $t = 1002$  to 0.014 in generation  $t = 1006$ . The decrease in accuracy from one generation to the next was similar for the two models.

For both intermediate marker densities, the accuracy of GWpEBV in generation  $t = 1002$  was lower compared to

**Table 1: Selection accuracy ( $r$ ) and regression of TBV on GWEBV and GWpEBV over five generations for marker density 8Ne/M, and accuracy differences when a polygenic effect is included**

Generation	Accuracy of selection		Regression of TBV on GWEBV	
	$r_{\text{GWEBV}}^{i)} \pm \text{s.e.}$	$\Delta_r^{ii)} \pm \text{s.e.}$	$b_{\text{GWEBV}}^{i)} \pm \text{s.e.}$	$\Delta_b^{ii)} \pm \text{s.e.}$
$t = 1002$	<b>0.875</b> 0.006	<b>0.008</b> 0.003	<b>0.926</b> 0.008	<b>0.058</b> 0.012
$t = 1003$	<b>0.861</b> 0.007	<b>0.007</b> 0.005	<b>0.917</b> 0.010	<b>0.068</b> 0.014
$t = 1004$	<b>0.857</b> 0.007	<b>0.007</b> 0.006	<b>0.916</b> 0.010	<b>0.074</b> 0.013
$t = 1005$	<b>0.852</b> 0.008	<b>0.011</b> 0.005	<b>0.912</b> 0.012	<b>0.080</b> 0.011
$t = 1006$	<b>0.842</b> 0.009	<b>0.014</b> 0.006	<b>0.906</b> 0.011	<b>0.079</b> 0.010

<sup>i)</sup>GWEBV represents the genome-wide estimated breeding value without polygenes

<sup>ii)</sup> $\Delta$  represents the difference between genome-wide estimated breeding values including a polygenic effect (GWpEBV) and GWEBV ( $\Delta_r = r_{\text{GWpEBV}} - r_{\text{GWEBV}}$  and  $\Delta_b = b_{\text{GWpEBV}} - b_{\text{GWEBV}}$ )

**Table 2: Selection accuracy (r) and regression of TBV on GWEBV and GWpEBV over five generations for marker density 4Ne/M and accuracy differences when a polygenic effect is included**

Generation	Accuracy of selection		Regression of TBV on GWEBV	
	$r_{\text{GWEBV}}^{i)} \pm \text{s.e}$	$\Delta_r^{ii)} \pm \text{s.e}$	$b_{\text{GWEBV}}^{i)} \pm \text{s.e}$	$\Delta_b^{ii)} \pm \text{s.e}$
$t = 1002$	<b>0.795</b> 0.006	<b>-0.014</b> 0.003	<b>0.896</b> 0.010	<b>-0.027</b> 0.009
$t = 1003$	<b>0.756</b> 0.006	<b>-0.002</b> 0.007	<b>0.864</b> 0.011	<b>0.032</b> 0.015
$t = 1004$	<b>0.732</b> 0.007	<b>0.006</b> 0.007	<b>0.848</b> 0.011	<b>0.065</b> 0.013
$t = 1005$	<b>0.722</b> 0.007	<b>0.011</b> 0.006	<b>0.846</b> 0.010	<b>0.075</b> 0.012
$t = 1006$	<b>0.705</b> 0.008	<b>0.014</b> 0.007	<b>0.827</b> 0.010	<b>0.095</b> 0.011

<sup>i)</sup>GWEBV represents the genome-wide estimated breeding value without polygenes

<sup>ii)</sup> $\Delta$  represents the difference between genome-wide estimated breeding values including a polygenic effect (GWpEBV) and GWEBV ( $\Delta_r = r_{\text{GWpEBV}} - r_{\text{GWEBV}}$  and  $\Delta_b = b_{\text{GWpEBV}} - b_{\text{GWEBV}}$ )

that of GWEBV ( $p < 0.05$ ), but after five generations of selection, the accuracy was significantly higher for GWpEBV (Tables 2 and 3) ( $p < 0.05$ ). For example, for marker density 4Ne/M, the difference in accuracy between GWpEBV and GWEBV was -0.014 in generation  $t = 1002$  and 0.014 after five generations (Table 2), indicating that by including a polygenic effect retained greater accuracy over generations. The accuracies for marker density 2 Ne/

M are relatively high compared to those for higher marker densities, which may be due to the structure of the marker/QTL map. At a marker density of 2 Ne/M, every SNP is adjacent to a putative QTL, whereas at higher densities a fraction of the SNP is not adjacent to any QTL [13].

For the lowest marker density (1Ne/M), the accuracy was 0.679 in the first generation for GWEBV and reduced to

**Table 3: Selection accuracy (r) and regression of TBV on GWEBV and GWpEBV over five generations for marker density 2Ne/M and accuracy differences when a polygenic effect is included**

Generation	Accuracy of selection		Regression of TBV on GWEBV	
	$r_{\text{GWEBV}}^{i)} \pm \text{s.e}$	$\Delta_r^{ii)} \pm \text{s.e}$	$b_{\text{GWEBV}}^{i)} \pm \text{s.e}$	$\Delta_b^{ii)} \pm \text{s.e}$
$t = 1002$	<b>0.801</b> 0.008	<b>-0.014</b> 0.006	<b>0.889</b> 0.011	<b>-0.074</b> 0.014
$t = 1003$	<b>0.763</b> 0.010	<b>0.007</b> 0.008	<b>0.862</b> 0.011	<b>-0.001</b> 0.012
$t = 1004$	<b>0.736</b> 0.011	<b>0.026</b> 0.006	<b>0.837</b> 0.013	<b>0.048</b> 0.014
$t = 1005$	<b>0.722</b> 0.011	<b>0.036</b> 0.008	<b>0.814</b> 0.012	<b>0.106</b> 0.009
$t = 1006$	<b>0.717</b> 0.009	<b>0.036</b> 0.007	<b>0.811</b> 0.010	<b>0.104</b> 0.010

<sup>i)</sup>GWEBV represents the genome-wide estimated breeding value without polygenes

<sup>ii)</sup> $\Delta$  represents the difference between genome-wide estimated breeding values including a polygenic effect (GWpEBV) and GWEBV ( $\Delta_r = r_{\text{GWpEBV}} - r_{\text{GWEBV}}$  and  $\Delta_b = b_{\text{GWpEBV}} - b_{\text{GWEBV}}$ )

**Table 4: Selection accuracy ( $r$ ) and regression of TBV on GWEBV and GWpEBV ( $b$ ) over five generations for marker density 1Ne/M and the accuracy differences when a polygenic effect is included**

Generation	Accuracy of selection		Regression of TBV on GWEBV	
	$r_{\text{GWEBV}}^{i)} \pm \text{s.e}$	$\Delta_r^{ii)} \pm \text{s.e}$	$b_{\text{GWEBV}}^{i)} \pm \text{s.e}$	$\Delta_b^{ii)} \pm \text{s.e}$
$t = 1002$	<b>0.679</b> 0.006	<b>0.005</b> 0.008	<b>0.866</b> 0.014	<b>-0.016</b> 0.013
$t = 1003$	<b>0.610</b> 0.009	<b>0.011</b> 0.012	<b>0.794</b> 0.017	<b>0.056</b> 0.026
$t = 1004$	<b>0.565</b> 0.013	<b>0.010</b> 0.013	<b>0.732</b> 0.015	<b>0.088</b> 0.023
$t = 1005$	<b>0.535</b> 0.013	<b>0.007</b> 0.015	<b>0.701</b> 0.016	<b>0.089</b> 0.024
$t = 1006$	<b>0.518</b> 0.012	<b>0.013</b> 0.013	<b>0.684</b> 0.016	<b>0.101</b> 0.025

<sup>i)</sup>GWEBV represents the genome-wide estimated breeding value without polygenes

<sup>ii)</sup> $\Delta$  represents the difference between genome-wide estimated breeding values including a polygenic effect (GWpEBV) and GWEBV ( $\Delta_r = r_{\text{GWpEBV}} - r_{\text{GWEBV}}$  and  $\Delta_b = b_{\text{GWpEBV}} - b_{\text{GWEBV}}$ )

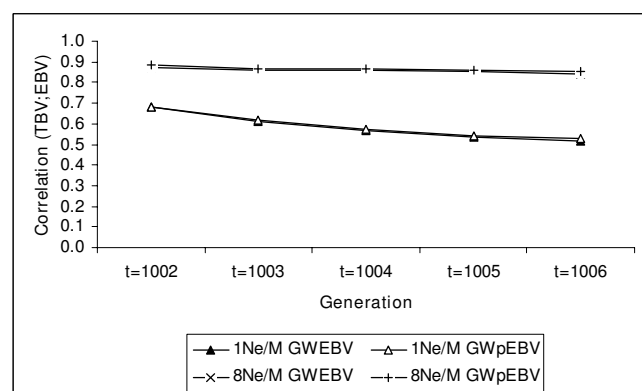
0.518 in generation  $t = 1006$  (Table 4). The use of GWpEBV increased the accuracy from 0.005 in generation  $t = 1002$  to 0.013 in generation  $t = 1006$ .

In general, when the polygenic effect was included in the model, the accuracy of GWpEBV was reduced in later generations as for GWEBV, but the decrease in accuracy was smaller for GWpEBV, especially for the intermediate marker densities. Figure 1 illustrates the selection accuracy over five generations for GWpEBV and GWEBV for the

highest and lowest marker densities, and the difference in accuracy between the two marker densities increased as the number of generations increased. Figure 1 clearly shows marginal differences between GWpEBV and GWEBV for the two marker densities, since the lines overlap, and that the accuracy is more stable over generations using a high marker density compared to a low marker density.

#### Regression coefficient of TBV on GWEBV and GWpEBV

The linear regression coefficient of TBV on GWEBV and GWpEBV was used as a measure of bias for these two selection criteria. Table 1, 2, 3, 4 show the regression coefficients of TBV on GWEBV and GWpEBV and the difference between the two models. For the highest marker density (8Ne/M), the regression coefficient for GWEBV was 0.926 in generation  $t = 1002$  and reduced to 0.902 in generation  $t = 1006$  (Table 1). The regression coefficient was significantly higher for GWpEBV than for GWEBV for all generations, and the difference between the models varied from 0.058 in generation  $t = 1002$  to 0.079 in generation  $t = 1006$ , respectively. Consequently the regression coefficients for GWpEBV were always between 0.98 and 1.00, i.e. showing only a very small bias. The reduction in regression was larger for GWEBV than for GWpEBV, as the regression coefficient for GWpEBV was much more stable over generations.



**Figure 1**  
**Accuracy of selection over five generations in different models.** Selection accuracy was determined for marker densities of 1Ne/M and 8Ne/M with the polygenic effect included (GWpEBV) or not (GWEBV) in the model; the lines for GWEBV and GWpEBV overlap almost completely, indicating minor differences between the two models.

For the intermediate marker densities, the regression coefficients were lower. However, there was a marked interaction between generation and method. For GWpEBV the regression coefficient was smaller than that for GWEBV at

generation  $t = 1002$ , but increased slightly over generations. In contrast, the regression coefficients for GWEBV decreased steadily over generations (Tables 2 and 3). By generation  $t = 1006$ , the difference in regression coefficients between GWEBV and GWpEBV was substantial: 0.095 (s.e. = 0.011) and 0.104 (s.e. = 0.010) for 4Ne/M and 2Ne/M, respectively. For 1Ne/M, both methods showed the same trend i.e. a decreasing regression coefficient, but the rate of decrease was faster for GWEBV.

In general, if the polygenic effect was ignored, bias increased from generation  $t = 1002$  to  $t = 1006$  for all four marker densities. However, this bias decreased with increasing marker densities (Table 1, 2, 3, 4). If a polygenic effect was included, the situation was similar, but the bias for all marker densities was more stable over generations, and furthermore decreased for intermediate generations (Table 2 and 3). For marker density 8Ne/M, the regression coefficient remained between 0.98 and 1.00 for all generations. Figure 2 shows the regression coefficient of TBV on GWEBV and GWpEBV for the highest marker density compared to the lowest marker density, and clearly shows that the regression coefficient is more stable over generations when the marker density is high and a polygenic effect is included.

### Polygenic variance

Table 5 shows the mean and median of the polygenic variances for the four different marker densities. The estimate of the mean polygenic variance ranged from 0.267 for 1Ne/M to 0.411 for 8 Ne/M with large standard errors. There was no statistically significant difference between the different marker densities, and no statistical evidence of a trend with increasing marker density. The distribution of the values for the Gibbs sampling within a replicate was very large with sporadic extreme values. This prompted us

**Table 5: Mean and median polygenic variance for the different marker densities in the base generation ( $t = 996$ ), estimated from the analysis of phenotypes in generation  $t = 1001$**

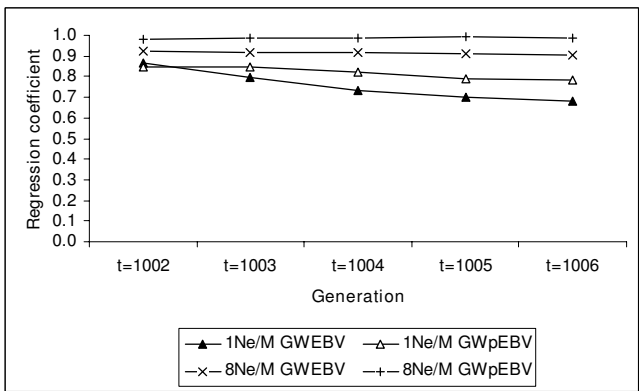
	1Ne/M	2Ne/M	4Ne/M	8Ne/M
Mean (s.e.)	<b>0.267</b> (0.090)	<b>0.323</b> (0.056)	<b>0.360</b> (0.028)	<b>0.411</b> (0.070)
Median (s.e.)	<b>0.266</b> (0.099)	<b>0.272</b> (0.059)	<b>0.252</b> (0.022)	<b>0.403</b> (0.082)

to examine the medians for the Gibbs samples for each replicate, they may be more robust to such outliers; however the picture changed very little.

### Discussion

This study shows that including a polygenic effect has little impact on the accuracy of genome-wide EBVs in the generation immediately following phenotyping. However, as the generations progress, the predictions with the polygenic effect retains somewhat greater accuracy. This persistence in accuracy over time is particularly significant for higher marker densities. This is because spurious marker associations arising from the pedigree are reduced, so that the remaining marker associations reflect more truly LD through proximity on the chromosome, which changes more slowly over time. Likewise the bias of the GWpEBV is significantly reduced compared to GWEBV, and the reduction is larger for the lowest marker densities, which displayed the largest bias for GWEBV. With lower marker densities, there are fewer markers around the QTL to explain the effect of the QTL, and the polygenic variance is expected to be more important for providing information for the estimated breeding values.

In Calus *et al.* [12], the accuracy of genomic selection including a polygenic effect was related to linkage disequilibrium (LD) between adjacent markers. For a high heritability trait, they found that including a polygenic effect increased selection accuracy when the  $r^2$  was lower than 0.14, and the benefit of including a polygenic effect increased with reduced  $r^2$ . The latter is consistent with our results in generation  $t = 1002$ , which is the only generation that can be compared to this study. In Calus *et al.* [12], the simulated model was based on a lower number of markers, smaller genome size and did not study the ability to predict GWEBV over multiple generations. The  $r^2$  values were calculated for a very similar dataset in an earlier paper, and were between 0.16 and 0.35 [8], which are larger than what Calus *et al.* [12] reported. As found here, the advantage of including a polygenic effect is more limited in the first generation after estimating marker effects. However, in practical situations, it may be advantageous to estimate the marker effects in one generation (e.g., due to phenotypic costs), and use these effects to



**Figure 2**  
**Regression coefficients of TBV on GWEBV and GWpEBV (bias) over five generations for marker densities 1Ne/M and 8Ne/M.**

select animals over multiple generations. Under such circumstances, it would be advantageous to include a polygenic effect since the accuracy will increase and bias decrease.

Whilst the accuracy of selection is a primary parameter of interest in animal breeding, the bias is also relevant since it determines the model's ability to predict the genetic progress. When generations are overlapping, individuals with different amounts of information and genetic level need to be compared for selection and biases, which can reduce the accuracies in predicting breeding values. Our results indicate that the polygenic effect did account for some of the variance not captured by the markers. Since estimates of polygenic effects are based on the BLUP theory and will thus show small bias, it may be expected and was found that including a polygenic effect reduces the bias.

The estimates of the variance of the polygenic effect increase with increasing marker density (Table 5), which is contrary to our expectation that, as marker density increases, the QTL will be more closely modelled by the markers and polygenic effects will become less important. A possible explanation for this result is that the non-linear regression implied by BayesB estimation of marker effects becomes more non-linear as marker density increases (because the fraction of markers with non-zero effect is expected to decrease). The increased non-linearity of the regression implies that small spurious associations will be increasingly regressed back to zero, resulting in more variance being explained by the polygenic effect. Furthermore, on a per marker basis, the spurious associations become smaller, since they are spread over more markers. These reductions in marker effects due to spurious associations may result in an increased variance attributed to the polygenic effect. This explanation implies that the effect of including or excluding a polygenic effect in these Bayesian models may depend on the prior distributions used for the marker effects and the polygenic effect, and different prior distributions may result in different outcomes. Also the number of QTL simulated (50-60) may affect the importance of the polygenic effect. It may be expected that with more QTL, the genetic model will become more like the infinitesimal model and the inclusion of a polygenic effect may be more beneficial.

Depending on the cost of genotyping and numbers of markers used, genomic selection programs will be more cost effective if the estimated marker effects could be used over multiple generations. Recombination will occur between the markers and QTL over time, resulting in reduced  $r^2$  and reduction in the accuracy of selection. This study shows that a marker density of 8Ne/M seems sufficient for the estimated marker effects to persist over five

generations with minimum bias and only a small reduction in selection accuracy. However, in practice the results will depend on the genetic architecture of the genome and on how similar the simulated parameters used in the study are compared to real genomes. Nevertheless, including a polygenic effect is beneficial for a random mating population when estimated marker effects are used to predict GWEBV over multiple generations, especially with respect to the bias.

## Competing interests

The authors declare that they have no competing interests.

## Authors' contributions

TRS simulated the datasets, carried out the analysis and drafted the manuscript. THEM wrote the computer modules and helped to carry out the study and draft the manuscript. All authors have read and approved the final manuscript.

## References

- Gianola D, Perez-Enciso M, Toro MA: **On marker-assisted prediction of genetic value: Beyond the ridge.** *Genetics* 2003, **163**:374-365.
- Gianola D, Fernando RL, Stella A: **Genomic-assisted prediction of genetic value with semiparametric procedures.** *Genetics* 2006, **173**:1761-1776.
- Meuwissen THE, Hayes BJ, Goddard ME: **Prediction of total genetic value using genome-wide dense marker maps.** *Genetics* 2001, **157**:1819-1829.
- Solberg TR, Sonesson AK, Woolliams JA, Meuwissen THE: **Reducing dimensionality for prediction of genome-wide breeding values.** *Genet Sel Evol* 2009, **41**:29.
- Calus MPL, Meuwissen THE, de Roos APW, Veerkamp RF: **Accuracy of genomic selection using different methods to define haplotypes.** *Genetics* 2008, **178**:553-561.
- Habier D, Fernando RL, Dekkers JCM: **The impact of genetic relationship information on genome-assisted breeding values.** *Genetics* 2007, **177**:2389-2397.
- Muir WM: **Comparison of genomic and traditional BLUP-estimated breeding value accuracy and selection response under alternative trait and genomic parameters.** *J Anim Breed Genet* 2007, **124**:342-355.
- Solberg TR, Sonesson AK, Woolliams JA, Meuwissen THE: **Genomic selection using different marker types and densities.** *J Anim Sci* 2008, **86**:2447-2454.
- Gjedrem T: *Selection and Breeding Programs in Aquaculture* Springer, Dordrecht, The Netherlands; 2005. 10-1-4020-3341-9
- Lynch M, Walsh B: *Genetics and Analysis of Quantitative Traits* Sinauer Associates, Inc. Massachusetts, USA; 1998.
- Sonesson AK, Meuwissen THE: **Testing strategies for genomic selection in aquaculture breeding programs.** *Genet Sel Evol* 2009, **41**:37.
- Calus MPL, Veerkamp RF: **Accuracy of breeding values when using and ignoring the polygenic effect in genomic breeding value estimation with a marker density of one SNP per cM.** *J Anim Breed Genet* 2007, **124**:362-368.
- Hayes BJ, Goddard ME: **The distribution of the effects of genes affecting quantitative traits in livestock.** *Genet Sel Evol* 2001, **33**:209-229.
- Gilks WR, Richardson S, Spiegelhalter DJ: *Markov chain Monte Carlo in Practice* Chapman & Hall/CRC, London, UK; 1996. 0-412-05551-1
- Sørensen D, Gianola D: *Likelihood, Bayesian, and MCMC Methods in Quantitative Genetics* Springer-Verlag, New York, USA; 2002. 0-387-95440-6
- Meuwissen THE, Luo Z: **Computing inbreeding coefficients in large populations.** *Genet Sel Evol* 1992, **24**:305-313.