



HAL
open science

Genetic evaluation using single-step genomic best linear unbiased predictor in American Angus

D.A.L Lourenco, Shogo Tsuruta, B.O. Fragomeni, Yutaka Masuda, Ignacio Aguilar, Andres Legarra, J.K. Bertrand, T.S. Amen, L. Wang, D.W. Moser, et al.

► To cite this version:

D.A.L Lourenco, Shogo Tsuruta, B.O. Fragomeni, Yutaka Masuda, Ignacio Aguilar, et al.. Genetic evaluation using single-step genomic best linear unbiased predictor in American Angus. *Journal of Animal Science*, 2015, 93 (6), pp.2653-2662. 10.2527/jas.2014-8836 . hal-02639641

HAL Id: hal-02639641

<https://hal.inrae.fr/hal-02639641>

Submitted on 28 May 2020

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial - NoDerivatives 4.0 International License

Running head: Single-step genomic BLUP for Angus**Genetic** evaluation using single-step genomic BLUP in American Angus¹

**D. A. L. Lourenco*², S. Tsuruta*, B. O. Fragomeni*, Y. Masuda*, I. Aguilar[†], A. Legarra[‡],
J. K. Bertrand*, T. S. Amen[§], L. Wang[§], D. W. Moser[§], and I. Misztal***

*Department of Animal and Dairy Science, University of Georgia, Athens, GA 30602

[†]Instituto Nacional de Investigacion Agropecuaria, Canelones, Uruguay 90200

[‡]Institut National de la Recherche Agronomique, UMR1388 GenPhySE, Castanet Tolosan,
France 31326

[§]Angus Genetics Inc., St. Joseph, MO 64506

¹This study was partially funded by the American Angus Association (St. Joseph, MO), Zoetis (Kalamazoo, MI), and by Agriculture and Food Research Initiative Competitive Grants no. 2015-67015-22936 from the US Department of Agriculture's National Institute of Food and Agriculture. **We gratefully acknowledge the very helpful comments by the two anonymous reviewers.**

²Corresponding author: danielino@uga.edu

1 **ABSTRACT:** Predictive ability of genomic EBV when using single-step genomic BLUP
2 (ssGBLUP) in Angus cattle was investigated. Over 6 million records were available on birth
3 weight (BW) and weaning weight (WW), almost 3.4 million on post-weaning gain (PWG), and
4 over 1.3 million on calving ease (CE). Genomic information was available on at most 51,883
5 animals, which included high and low EBV accuracy animals. Traditional EBV was computed
6 by BLUP and genomic EBV by ssGBLUP and indirect prediction based on SNP effects derived
7 from ssGBLUP; SNP effects were calculated based on the following reference populations:
8 ref_2k (high EBV accuracy sires and cows), ref_8k (ref_2k, plus all genotyped ancestors of
9 validation animals), and ref_33k (ref_8k, plus all remaining genotyped animals not in the
10 validation). Indirect prediction was obtained as direct genomic value (DGV) or as an index of
11 DGV and parent average (PA). Additionally, runs with ssGBLUP used the inverse of the
12 genomic relationship matrix calculated by an algorithm for proven and young animals (APY)
13 that uses recursions on a small subset of reference animals. An extra reference subset included
14 3872 genotyped parents of genotyped animals (ref_4k). Cross-validation was used to assess
15 predictive ability on a validation population of 18,721 animals born in 2013. Computations for
16 growth traits used multiple-trait linear model, and for CE, a bivariate CE-BW threshold-linear
17 model. With BLUP, predictivities were 0.29, 0.34, 0.23, and 0.12 for BW, WW, PWG, and CE,
18 respectively. With ssGBLUP and ref_2k (ref_33k), predictivities were 0.34, 0.35, 0.27, and 0.13
19 (0.39, 0.38, 0.29, and 0.13), respectively. Low predictivity for CE was due to low incidence rate
20 of difficult calving. Indirect predictions with ref_33k were as accurate as with full ssGBLUP.
21 Using APY and recursions on ref_4k (ref_8k) gave 88% (97%) gains of full ssGBLUP.
22 Genomic evaluation in beef cattle with ssGBLUP is feasible while keeping the models (maternal,
23 multiple trait, threshold) already used in regular BLUP. Gains in predictivity are dependent on

24 the **composition** of the reference population. Indirect predictions via SNP effects derived from
25 ssGBLUP allow for accurate genomic predictions on young animals, with no advantage of
26 including PA in the index if the reference population is large. With APY conditioning on **about**
27 10,000 reference animals, ssGBLUP is potentially applicable to large number of genotyped
28 animals without compromising **predictive ability**.

29 **Key words:** beef cattle, genomic recursion, genomic selection, indirect prediction

30 INTRODUCTION

31 Genomic selection in beef cattle has currently been performed with multistep methods, which
32 uses deregressed EBV to estimate **SNP effects and then** direct genomic value (DGV) for
33 selection candidates based on their genotypes (Meuwissen et al., 2001; Garrick et al., 2009). The
34 main advantage of this approach is that the traditional BLUP evaluation is kept unchanged and
35 genomic selection can be carried out by a separate entity owning genotypes but not phenotypes.
36 Also new animals are easily evaluated if DGV is computed as a sum of marker effects, but not if
37 selection indexes including DGV and parent average (PA) are used.

38 When both phenotypes and genotypes are available jointly, single-step genomic BLUP
39 (ssGBLUP) (Aguilar et al., 2010) **is a simple alternative**. This method does not rely on
40 deregressed proofs, properly weighs information from genotyped sires and cows, **thus avoiding**
41 **double-counting of contributions due to relationships and records**, and accounts for pre-selection
42 bias of genomically selected parents without phenotypes (Legarra et al., 2014). **In ssGBLUP it is**
43 **also possible to quickly evaluate young genotyped animals without running a complete**
44 **evaluation that requires several hours to converge. Quick predictions can be calculated indirectly,**

45 where genomic predictions for young animals are obtained from SNP effects. It was shown by
46 Wang et al. (2012) that SNP effects can be derived from GEBV solutions from the main
47 ssGBLUP evaluation.

48 In its current implementation, ssGBLUP uses direct inversion of genomic matrices (Aguilar
49 et al., 2011), which has a cubic cost and a limit of 150,000 animals (Aguilar et al., 2013). Several
50 methods were proposed to overcome that limit (Legarra and Ducrocq, 2012; Fernando et al.,
51 2014; Liu et al., 2014), but none was successful. Recently Misztal et al. (2014) presented a
52 method which uses an approximate inversion of genomic relationships based on recursions on a
53 fraction of the total population; which can be suitable and inexpensive. The first goal of this
54 study was to evaluate the feasibility of ssGBLUP for genomic evaluation in Angus cattle with
55 reference populations of different composition. An additional goal was to evaluate the ability to
56 predictive GEBV with genomic recursions and with indirect prediction for young animals.

57 MATERIAL AND METHODS

58 Datasets from American Angus Association (AAA) were available for this study that
59 included growth traits and calving ease (CE). Growth traits included birth weight (BW), weaning
60 weight (WW), and post-weaning gain (PWG). As the data were obtained from existing
61 databases, Animal Care and Use Committee approval was not obtained for this study.

62 *Data*

63 Over 6 million phenotypes were available for BW and WW, almost 3.4 million for PWG, and
64 over 1.3 million for CE. Whereas BW, WW, and PWG are continuous traits, CE is a categorical
65 trait with 5 calving scores, where 5 is abnormal delivery and is excluded. Because few animals
66 had scores 3 and 4, these scores were combined into category 2, which resulted in 93% of

67 animals with score 1 and 7% with score 2. The number of animals in the pedigree for evaluation
68 of growth traits was 8,236,425, and for CE was 8,025,676.

69 For evaluation of growth traits, 81,878 animals were genotyped for 54,609 SNP from the
70 BovineSNP50k v2 BeadChip (Illumina Inc., San Diego, CA). Currently, no genotyping strategy
71 is applied by AAA; therefore, the members can choose which animals are being genotyped, and
72 most of them are young. A total of 29,995 genotyped animals were young without phenotypes
73 for any of the 3 traits, which caused them to have genotypes excluded from this study. If the
74 number of genotyped animals is relatively large, young genotyped animals without phenotypes
75 in the dataset give very small contribution to their relatives' evaluation (Misztal et al., 2014).
76 After removing SNP with unknown position or located on sex chromosomes and running a
77 general quality control analysis, genotypes on 38,528 SNP markers were available for 32,465
78 males and 19,418 females born from 1977 to 2013; therefore, the maximum number of
79 genotyped animals used in all analyses on growth traits was 51,883. For CE evaluation, a
80 genotyping set with 72,069 animals was available, but only genotypes on 40,546 animals born
81 from 1977 to 2013 (26,074 males and 14,472 females) were used for the same reason above. The
82 number of SNP that passed the general quality control for this dataset was 38,568.

83 For this study, the animals were then split into training and validation populations according
84 to year of birth. Thus, all 18,721 (13,166) genotyped animals born in 2013 were chosen to be in
85 the validation population for growth (CE) traits and had their phenotypes removed from the
86 evaluations. The pedigree relationship between training and validation populations ranged from 0
87 to 0.82, with an average relationship of 0.09.

88 ***Model***

89 Traditional and genomic evaluations were performed for growth traits and CE. A
90 multivariate linear animal model was used for growth traits as:

$$91 \quad \mathbf{y}_t = \mathbf{X}\mathbf{b} + \mathbf{Z}_1\mathbf{u} + \mathbf{Z}_2\mathbf{m} + \mathbf{Z}_3\mathbf{p} + \mathbf{e} \quad [1]$$

92 where t is for each one of BW, WW, PWG; \mathbf{y} , \mathbf{b} , \mathbf{u} , \mathbf{m} , \mathbf{p} , and \mathbf{e} are vectors of phenotypes, fixed
93 effect of contemporary group, additive direct genetic effect, additive maternal genetic effect,
94 maternal permanent environmental effect, and random residuals, respectively; \mathbf{X} , \mathbf{Z}_1 , \mathbf{Z}_2 , and \mathbf{Z}_3
95 are incidence matrices for \mathbf{b} , \mathbf{u} , \mathbf{m} , and \mathbf{p} , respectively. All random effects were present for WW,
96 but only \mathbf{u} , \mathbf{m} , and \mathbf{e} for BW, and \mathbf{u} and \mathbf{e} for PWG.

97 A bivariate threshold-linear animal model was used to model CE jointly with BW:

$$98 \quad \mathbf{Y}_c = \mathbf{X}\mathbf{b} + \mathbf{Z}_1\mathbf{u} + \mathbf{Z}_2\mathbf{m} + \mathbf{e} \quad [2]$$

99 where c is for BW and CE; \mathbf{y} , \mathbf{b} , \mathbf{u} , \mathbf{m} and \mathbf{e} are vectors of phenotypes, fixed effects of
100 contemporary group, sex, age of dam (only for CE), and sex by age of dam interaction (only for
101 CE), additive direct genetic effect, additive maternal genetic effect, and random residuals,
102 respectively; \mathbf{X} , \mathbf{Z}_1 , and \mathbf{Z}_2 are incidence matrices for \mathbf{b} , \mathbf{u} , and \mathbf{m} , respectively. According to
103 Ramirez-Valverde et al. (2001) when BW is available, bivariate threshold-linear models
104 including CE and BW are a better alternative than a single-trait threshold model to evaluate CE,
105 especially if the population has animals with different levels of EBV accuracy. From this model,
106 only results for CE are discussed, whereas results for BW are from the multiple trait linear model
107 for growth traits. Heritabilities for all traits were calculated by AAA using the same models as in
108 [1] for BW, WW, and PWG; and in [2] for CE. For our study, the values were then provided by
109 AAA and ranged from 0.12 to 0.41 (Table 1).

110 *Analyses*

111 Three different genomic analyses were performed using ssGBLUP (Aguilar et al., 2010;
112 Christensen and Lund, 2010) as implemented in BLUP90IOD program
113 (<http://nce.ads.uga.edu/wiki/BLUPmanual>). Compared to BLUP, in ssGBLUP the inverse of the
114 numerator relationship matrix \mathbf{A}^{-1} is replaced by matrix \mathbf{H}^{-1} defined as follows:

$$115 \quad \mathbf{H}^{-1} = \mathbf{A}^{-1} + \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{G}^{-1} - \mathbf{A}_{22}^{-1} \end{bmatrix},$$

116 where \mathbf{G} is the genomic relationship matrix. The computations used default options in
117 BLUP90IOD. In all analyses the validation population was defined as genotyped animals born in
118 2013 with phenotypes excluded.

119 ***First analysis: ssGBLUP with different reference populations.*** Different reference
120 populations were defined according to EBV accuracy calculated with the ACCF90 program
121 (<http://nce.ads.uga.edu/wiki/BLUPmanual>), which uses the concept of prediction error variance
122 and reflects the standard error of EBV for each individual. The objective was to investigate the
123 influence of different groups of reference animals on genomic predictions, and possibly to guide
124 genotyping strategy. The current trend in livestock genomics is to genotype young animals;
125 however, more important animals give more information to the evaluations. For growth traits
126 (CE), the first reference population was composed of 1,628 (1,541) top bulls with EBV accuracy
127 for $BW \geq 0.85$; which we will refer hereinafter as "ref_bulls". As BW was present in models for
128 growth and CE evaluations, using its EBV accuracy for selecting top bulls helped to compose
129 sets with proportional number of animals. In this case, the \mathbf{G} matrix was composed of animals in
130 the reference population and also animals in the validation population; the last had 18,721
131 animals for growth traits and 13,166 for CE. The second reference population was composed of
132 the top bulls and also top cows that had an EBV accuracy for $BW \geq 0.85$; which we will refer as

133 **ref_2k**. The number of top cows was small and only 268 were added for the growth trait analysis
134 and 323 for CE. The third reference population was composed of top bulls, top cows, and all
135 other genotyped animals born from 1977 to 2012 (**we will refer as ref_33k**). This group had a
136 total of 33,162 animals for growth and 27,380 for CE, with an average **EBV accuracy for BW** of
137 0.77 (± 0.16). For the latter analysis, the **G** matrix was composed of the maximum number of
138 51,883 genotyped animals for growth analysis and 40,546 for analysis of CE.

139 ***Second analysis: ssGBLUP with indirect predictions for young animals.*** With the
140 increasing number of genotyped heifers and steers in dairy and beef, the genomic methods
141 should be able to provide predictions for young animals without phenotypes in a quick run,
142 externally to the official evaluations. This concept is introduced here as indirect ssGBLUP, and
143 basically mimics the mixed model equations. **It would be advantageous from different**
144 **perspectives: to evaluate young animals mainly for traits that are measured later in life, after the**
145 **selection decisions are made; and to reduce computing costs because the dimension of **G** would**
146 **not increase in the same proportion as the number of genotyped animals.**

147 In order to explain how it works, consider the equation for the GEBV of a single individual
148 in ssGBLUP (VanRaden and Wiggans, 1991; Aguilar et al., 2010; Lourenco et al., 2015):

$$149 \text{GEBV} = w_1\text{PA} + w_2\text{YD} + w_3\text{PC} + w_4\text{DGV} - w_5\text{PP}$$

150 where PA is Parent Average, YD is Yield Deviation (phenotypes adjusted for model effects other
151 than additive genetic and error), PC is Progeny Contribution, DGV is direct genomic value
152 (corresponding to \mathbf{G}^{-1}), PP is the pedigree prediction based on the subset of genotyped animals
153 from **A** (corresponding to \mathbf{A}_{22}^{-1}) and w_1 to w_5 are weights that sum to 1. In the case of young
154 animals with no progeny or own performance record, $\text{YD}=\text{PC}=0$ and $w_2=w_3=0$. In this case, for

155 individual i , $PA_i = (GEBV_s + GEBV_d)/2$; $DGV_i = -\frac{\sum_{j,j \neq i} g^{ij} GEBV^j}{g^{ii}}$; $PP = -\frac{\sum_{j,j \neq i} a_{22}^{ij} GEBV^j}{a_{22}^{ii}}$ and $w_l = \frac{2}{den}$,

156 $w_4 = \frac{g^{ii}}{den}$, $w_5 = \frac{a_{22}^{ii}}{den}$, where den is the denominator that equals to $2 + (g^{ii} - a_{22}^{ii})$; g^{ij} (a_{22}^{ij}) is an

157 element of \mathbf{G}^{-1} (\mathbf{A}_{22}^{-1}) corresponding to relationships between animal i and j ; s and d correspond to

158 sire and dam, respectively. If all individuals are genotyped, then $PA=PP$ and GEBV reduces to

159 DGV.

160 For ssGBLUP with indirect predictions, SNP effects can be calculated using the current run

161 of ssGBLUP with all but young animals, and genomic predictions for young animals are

162 obtained by multiplying the SNP content by SNP effect to obtain DGV; a more complete GEBV

163 can also be available through a selection index that combines DGV and PA. The flow for indirect

164 predictions in ssGBLUP is:

165 1) Run ssGBLUP with a reference population to obtain GEBV. In this step, 3 reference
166 populations were tested:

167 a) ref_2k: reference population with top bulls and top cows ($n=1,896$);

168 b) ref_8k: reference population with all parents that were genotyped ($n=8,285$), this
169 includes ref_2k;

170 c) ref_33k: reference population with all genotyped animals born up to 2012 ($n=33,162$),
171 this includes ref_8k;

172 2) Split GEBV into all the components shown before, where DGV for an animal i in the

173 reference population is calculated as below (Aguilar et al. (2010):

$$DGV_i = -\frac{\sum_{j,j \neq i} g^{ij} GEBV^j}{g^{ii}}$$

174 with all elements previously defined.

175 3) Calculate SNP effects using DGV from the reference population:

$$\hat{\mathbf{u}} = \mathbf{DZ}'\mathbf{G}^{-1}(\mathbf{DGV})$$

176 where $\hat{\mathbf{u}}$ is a vector of estimated SNP effects, \mathbf{D} is a diagonal matrix of weights
 177 (standardized variances) for SNP (identity matrix in this case), and \mathbf{Z} is a matrix of
 178 centered genotypes for each animal (VanRaden, 2008). A similar approach that uses
 179 GEBV instead of DGV to calculate SNP effects was proposed by Wang et al. (2012).
 180 However, for numerical purposes this involves approximations as \mathbf{G} matrix is formed as
 181 $\mathbf{G} = 0.95\mathbf{ZDZ}' + 0.05\mathbf{A}_{22}$ (Aguilar et al., 2010). This is done as a default approach to avoid
 182 singularity problems and may result in negligible error as shown later.

183 4) Calculate DGV for young genotyped animals (\mathbf{DGV}_y):

$$\mathbf{DGV}_y = \mathbf{Z}_y \hat{\mathbf{u}}$$

184 where \mathbf{DGV}_y and \mathbf{Z}_y are direct genomic values and a matrix of centered genotypes for
 185 young animals not included in ssGBLUP evaluation, respectively.

186 5) Combine \mathbf{DGV}_y with PA for young genotyped animals:

$$\mathbf{GEBV}_y \approx w_1 \mathbf{PA} + w_4 \mathbf{DGV}_y$$

187 where \mathbf{GEBV}_y is GEBV obtained via indirect predictions for young animals, w_1 and w_4
 188 are weights identical for all animals and calculated based on covariances between \mathbf{DGV}_y
 189 and PA as:

$$\begin{bmatrix} w_1 \\ w_4 \end{bmatrix} = \begin{bmatrix} \sigma_{\text{PA}}^2 & \sigma_{\text{PA}, \text{DGV}_y} \\ \sigma_{\text{DGV}_y, \text{PA}} & \sigma_{\text{DGV}_y}^2 \end{bmatrix}^{-1} \begin{bmatrix} \sigma_{\text{PA}}^2 \\ \sigma_{\text{DGV}_y}^2 \end{bmatrix}$$

191 Note this is an approximation which ignores PP. In general, PP includes part of PA
 192 explained by DGV. When all animals are genotyped, PP and PA cancel out, with
 193 approximate cancellation when parents of an animal are genotyped. When an animal is
 194 unrelated to a genotyped population, PP=0. Fixed weights in the index account for an

195 average relationship of all young animals to a genotyped population. It is possible to
 196 create different indices based on the number of genotyped parents (VanRaden et al.,
 197 2012).

198 The ssGBLUP with indirect prediction allows calculation of DGV or GEBV for young
 199 genotyped animals, with lower computing cost compared to a full ssGBLUP where young
 200 animals are explicitly included.

201 ***Third analysis: ssGBLUP with G inverted by a recursive algorithm.*** When the number of
 202 genotyped animals is large and there is a desire for using all of them in ssGBLUP evaluations to
 203 get direct predictions for all, including young animals, an algorithm that splits genotypes into
 204 proven and young animals and uses recursion to approximate the inverse of the \mathbf{G} matrix was
 205 proposed by Misztal et al. (2014). This algorithm is known as APY, and \mathbf{G}^{-1} containing all
 206 genotyped animals can be expressed as:

$$\mathbf{G}^{-1} = \begin{bmatrix} \mathbf{G}_{pp}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} + \begin{bmatrix} -\mathbf{G}_{pp}^{-1} \mathbf{G}_{py} \\ \mathbf{I} \end{bmatrix} \mathbf{M}_g^{-1} \begin{bmatrix} -\mathbf{G}_{yp} \mathbf{G}_{pp}^{-1} & \mathbf{I} \end{bmatrix}$$

207 where the subscript pp stands for proven animals and py for the covariance between proven and
 208 young animals; each element of \mathbf{M}_g is obtained (for the i^{th} young animal) as

209 $m_{g,i} = g_{ii} - \mathbf{G}_{ip} \mathbf{G}_{pp}^{-1} \mathbf{G}_{pi}$ and is called genomic Mendelian sampling. In APY, the only direct
 210 inversion needed is for part of \mathbf{G} that contains relationships among proven animals (\mathbf{G}_{pp}),
 211 whereas all other coefficients are obtained through recursions.

212 For this analysis, four definitions of proven animals were tested that included the 3
 213 definitions used for indirect predictions (ref_2k, ref_8k, and ref_33k), plus one more definition
 214 where 3,872 genotyped parents of genotyped animals were considered as proven (ref_4k). This

215 last group was added to test if proven animals would have strong links with the young genotyped
216 population.

217 The greatest advantages of this algorithm are the reduction of computing cost, which is still
218 cubic for proven animals, but can be linear for young animals; and the possibility of using large
219 amounts of genotyped animals in ssGBLUP evaluations. The secondary advantage is numerical
220 stability as the regular \mathbf{G} matrix is singular when the number of animals is greater than the
221 number of SNP markers and cannot be inverted without blending with \mathbf{A}_{22} .

222 *Validation*

223 The ability to predict future phenotypes was the validation method chosen for this study. This
224 method is based on Legarra et al. (2008), and predictive ability for traditional and genomic
225 evaluations for animals born in 2013 was calculated as the correlation between (G)EBV and
226 phenotypes corrected for fixed effects ($y-Xb$):

$$r = \text{cor}[(G)EBV, y-Xb]$$

227 The predictive ability or predictivity is used as an approach to compare the methods applied
228 in this paper. For all analyses, the validation groups were kept the same to make comparisons
229 easier. Validations involved 18,721 animals for growth traits and 16,133 animals for CE.
230 Predictivity calculated with EBV in the above formula was the benchmark used to compare the
231 gain in predictive ability due to genomics, and predictivity calculated with GEBV was used to
232 compare the genomic methods previously described. Prediction accuracy could be described as
233 r/h , where h is square root of heritability; however, prediction accuracy can be overestimated if
234 heritabilities are obtained by simplified models as the ones used by AAA.

235

RESULTS AND DISCUSSION

236 *ssGBLUP with different reference populations*

237 **Predictive ability** on young animal when using several reference populations is shown in
238 Table 2. Using only top bulls as a reference population (**ref_bulls**) increased predictivity relative
239 to BLUP **by 0.05 for BW, 0.01 for WW, 0.04 for PWG, and 0.01 for CE**. Addition of top cows to
240 the reference population (**ref_2k**) did not increase the **predictivity** for any trait. This could be due
241 to the small number of animals added and also because daughters of those cows already
242 contributed through the inclusion of bulls. Addition of around **31,000** animals to the reference
243 population provided an additional increase in **predictivity of 0.05 for BW, of 0.03 for WW and of**
244 **0.02 for PWG**. However, no additional increase was observed for CE by adding extra **27,000**
245 genotyped animals, of which about **7,000** had phenotypes for that trait.

246 The addition of **31,000** animals with few or no progeny led to the same increase of
247 **predictivity** as using only the top bulls for BW, led to an increase of 3 times for WW and an
248 increase of 0.5 times for PWG. Among the 31,000 extra animals, almost all had phenotypes for
249 BW and WW, but only 24,000 had phenotypes for PWG. **Evidently, the composition of reference**
250 **population is also a factor that influences predictivity of GEBV besides the reference population**
251 **size. Thus, genotyping strategy should take into account genotyping more important and maybe**
252 **older animals with more information (higher EBV accuracy) along with genotyping large**
253 **amounts of young animals.**

254 Previous studies showed that **prediction accuracies or predictive ability** are biased downward
255 by selection (Bijma, 2012; Lourenco et al., 2015). In our study, it appears that selection for
256 proven bulls was much stronger for WW than for PWG (**lower increase in predictivity** with twice
257 the phenotypic data at similar heritability) but there was a small selection on genotyped animals

258 with own records (approximately twice the increase of **predictivity** with twice the phenotypic
259 data). It may be hard to calculate the amount of bias in livestock species, including beef cattle, as
260 the selection process is sequential and affected by both genetic correlations and specific indexes
261 used for selection.

262 Low **predictivity** for CE in this study is due to lower heritability combined with limited
263 recording for this trait and a low incidence of difficult calving. Additionally, very few genotyped
264 animals had a difficult calving, **perhaps because animals from a difficult calving are unlikely to**
265 **be retained for breeding and therefore would not be genotyped on a regular basis.** Higher
266 **predictivity** and impact of genomic selection for CE could be expected in breeds with higher
267 incidence of calving problems.

268 **Because the increase in predictivity for CE was very small compared to predictivity of**
269 **traditional evaluations, indirect predictions and APY were not tested for this trait.**

270 **In this paper, only predictivity for the direct genetic effect is shown; however, models for**
271 **BW and WW included maternal effect, which is also important in genetic evaluations. We**
272 **attempted to derive formulas for predictivity of maternal effects, unsuccessfully. Such**
273 **predictivity can be hard to assess because the maternal effect occurs one generation back, which**
274 **means that the corrected phenotype of animal i should be correlated with the maternal effect of**
275 **the dam of animal i. But, dams usually have more than one progeny and there is genetic**
276 **correlation between direct and maternal for BW, which makes derivations difficult. Lourenco et**
277 **al. (2013) used simulated data that mimicked a beef cattle population and showed that the gain**
278 **for the maternal effect with ssGBLUP is as high as for the direct effect.**

279 *ssGBLUP with indirect predictions for young animals*

280 Predictive ability for indirect prediction via conversion of DGV into SNP effects is shown in
281 Figure 1. When the reference population included top bulls and top cows (ref_2k), the
282 predictivity of indirect DGV_y was lower than predictivity for traditional EBV for the three traits
283 (0.23 vs. 0.29 for BW; 0.28 vs. 0.34 for WW; 0.19 vs. 0.23 for PWG). Predictivity for $GEBV_y$
284 calculated as an index of indirect DGV_y with PA was higher than those for EBV for the three
285 traits (0.31 vs. 0.29 for BW; 0.36 vs. 0.34 for WW; 0.24 vs. 0.23 for PWG), however, this
286 predictivity was lower than the ones from full ssGBLUP (except for WW). With larger reference
287 population (ref_8k), all indirect DGV_y were similar or more accurate than EBV, and the index
288 had similar predictivity as the full ssGBLUP. With the largest reference population (ref_33k), all
289 indirect DGV_y were almost as accurate as GEBV from full ssGBLUP, with the index marginally
290 improving predictivity for WW. This marginal improvement for WW may be caused by the use
291 of less than optimal genetic parameters, e.g., zero covariance between direct and maternal effects
292 (to reduce computing costs). The DGV_y obtained with ref_33k reference population were more
293 accurate than GEBV from full ssGBLUP obtained with ref_8k reference population.

294 Although predictivity of indirect predictions when using ref_33k was similar to predictivity
295 from full ssGBLUP, it does not mean that predictions have the same average. The reason for that
296 is the different sources of information used to calculate indirect predictions. Correlations
297 between GEBV and indirect predictions are a good tool to assure that the latter can be used for
298 interim evaluations. Correlations between GEBV from full ssGBLUP and DGV_y or $GEBV_y$ from
299 indirect predictions are shown in Table 3. On average, correlations with DGV_y were 0.73, 0.89,
300 and 0.96 for ref_2k, ref_8k, and ref_33k, respectively. Higher correlations were observed
301 between GEBV and $GEBV_y$, with values for the three reference sets being 0.89, 0.95, and 0.97,
302 respectively. Those results endorse the use of a reference population of size close to 33,000

303 animals for this American Angus dataset. By doing that, indirect predictions are as accurate as
304 predictions including genotypes for young animals in the evaluation (full ssGBLUP).

305 For young animals, $GEBV = w_1PA + w_4DGV - w_5PP$, with all weights summing to 1.0
306 (VanRaden and Wiggans, 1991; VanRaden et al., 2009; Aguilar et al., 2010). When the number
307 of genotyped animals is small, w_4 is small and ignoring PA reduces predictivity. Using an index
308 with PA improves the predictivity, however, PP is ignored and computed weights w_1 and w_4 are
309 approximate. When the number of genotyped animals is large, w_4 is close to 1.0, and ignoring
310 PA marginally reduces the predictivity for some traits. Therefore, the indirect prediction via
311 DGV is accurate when SNP effects are derived from ssGBLUP with sufficient size of the
312 reference population.

313 Neglecting PP seems to have no considerable effect in this population, because predictivity
314 of indirect predictions was very similar to predictivity from full ssGBLUP. Neglection of PP
315 indirectly means adjusting PA for an average PP. VanRaden et al. (2012) used different weights
316 for animals based on the number of genotyped parents, which better accounts for PP.

317 A study by Wiggans et al. (2014) used SNP effects from previous monthly genomic multi-
318 step evaluations to calculate preliminary GEBV for young genotyped animals. The objective was
319 to have daily or weekly genomic evaluations for US dairy cattle and reduce the time between
320 having DNA samples and predictions from a monthly official evaluation. Their reference set
321 contained all genotyped animals with phenotypes (about 597,000; corresponding to ref_33k in
322 our study) and correlations between preliminary and official evaluations were higher than 0.99
323 for Holsteins, but smaller for other breeds with a smaller number of genotyped animals. Further
324 research with different species will be critical in determining the sufficient size of the reference
325 population for indirect predictions in order to achieve high predictivity. It may be related to

326 effective population size, number of independent SNP (Pintus et al., 2013), and relationships
327 between reference and validation populations as in multi-step methods. Although indirect
328 predictions via ssGBLUP use a concept similar to multi-step methods for young genotyped
329 animals, indirect predictions via ssGBLUP may be more accurate than multistep predictions
330 because the latter are affected by approximations involved in deregressions and possible double-
331 counting of phenotypic information.

332 For young animals, indirect predictions via SNP effects from ssGBLUP seems a viable
333 alternative as it can be done separately from the full evaluation. As SNP effects are calculated
334 based on trait GEBV or DGV, indirect predictions are easily obtained for multi-trait models, as
335 done in this study; multi-breed and crossbred evaluations are possible when the G matrix is able
336 to account for information on all breeds. However, if young animals and particularly full-sibs are
337 intensively selected, selection on the Mendelian sampling will not be accounted for, leading to
338 pre-selection bias (Patry and Ducrocq, 2011). Analyses by ssGBLUP with all genotypes subject
339 to selection are expected to account for pre-selection (VanRaden and Wright, 2013), because
340 selection is accounted for when all information used for selection is included in the model
341 (Henderson, 1975).

342 *Comments on SNP weighting and SNP selection*

343 The way SNP effects are calculated in ssGBLUP allows for inclusion of different weights for
344 SNP: $\hat{\mathbf{u}} = \mathbf{DZ}'\mathbf{G}^{-1}(\mathbf{DGV})$, with weights for \mathbf{G} fit into the diagonal matrix \mathbf{D} . Those weights can be
345 calculated through an iterative process, or external weights can be used as input for ssGBLUP
346 (Wang et al., 2012; Su et al., 2014). Weighting \mathbf{G} seems to be a reasonable approach to achieve
347 higher prediction accuracy, especially in the presence of “major” SNP. Sun et al. (2011) showed

348 higher prediction accuracy when using weighted \mathbf{G} in regular GBLUP compared to BayesB. For
349 some traits, SNP weighting or SNP selection in ssGBLUP also gave additional **prediction**
350 **accuracy** (Wang et al., 2014). **In fact**, when weights are different per trait, this precludes the use
351 of multiple traits unless the model includes one common additive effect and specific additive
352 effects for individual traits. In practice and especially under a selection index, gains from a
353 multiple trait analysis can overcome losses due to not fitting “major” SNP. Also, when the
354 number of genotyped animals increases, the rate of gain in reliability increases at a slower pace
355 (VanRaden et al., 2011); therefore, weighting SNP may no longer have a big impact on
356 **prediction accuracy** (Winkelman et al., 2015).

357 *ssGBLUP with \mathbf{G} inverted by a recursive algorithm*

358 **Predictive ability** of GEBV when the inverse of \mathbf{G} is computed with APY is shown in Figure
359 2. When the recursions were conditioned on **ref_2k, ref_4k, ref_8k, and ref_33k**, the procedure
360 accounted for **67%, 88%, 97%, and 100%** of **predictivity** gains of ssGBLUP over BLUP,
361 respectively. Therefore, in ssGBLUP, using genomic recursion to invert \mathbf{G} while conditioning on
362 enough number of animals, in this case about 8,000, has the same prediction power as \mathbf{G} using
363 direct inversion. **The amount of memory necessary for APY \mathbf{G}^{-1} using ref_2k, ref_4k, ref_8k,**
364 **and ref_33k was approximately 0.8, 1.6, 3.2, and 13.7 Gbytes, respectively, whereas the amount**
365 **of memory for the regular \mathbf{G}^{-1} is 21.6 Gbytes. Therefore, using APY \mathbf{G}^{-1} makes computations**
366 **less costly and faster.**

367 Tests involving **100,000** genotyped Holsteins with recursions conditioned on more than
368 **15,000** animals resulted in practically identical GEBV compared to the regular inversion but with
369 a better convergence rate (Fragomeni et al., 2015) indicating that APY has good predictive and

370 numerical properties. They suggested that the necessary number of animals being conditioned is
371 proportional to the number of independent chromosome segments, which is a function of an
372 effective population size.

373 The main advantages of APY are low computing costs and numerical stability. With
374 conditioning on 8,000 animals, for example, the only inverse required is for a block of \mathbf{G} for
375 8,000 animals, and additional genotypes require only linear storage and computations.
376 Subsequently, computations with a large number of genotyped animals may be feasible with
377 similar predictivity as in the regular inversion. APY would be the algorithm of choice for regular
378 evaluations with very large number of genotyped animals.

379 CONCLUSIONS

380 Genomic evaluation in beef cattle using single-step genomic BLUP is feasible for either
381 linear or categorical traits. Gains in predictive ability over BLUP are dependent on the size and
382 composition of the reference population, and are large for growth traits and small for CE. With a
383 sufficient number of animals in the reference population, indirect prediction for young animals
384 via SNP effects provides similar predictivity to full single-step genomic BLUP, allowing for
385 quick genomic predictions without running a complete evaluation. Use of the algorithm for
386 proven and young animals in single-step genomic BLUP allows for incorporation of large
387 number of genotyped animals at low cost without compromising the predictive ability.

388 LITERATURE CITED

- 389 Aguilar, I., I. Misztal, D. L. Johnson, A. Legarra, S. Tsuruta, and T. J. Lawlor. 2010. Hot topic:
390 A unified approach to utilize phenotypic, full pedigree, and genomic information for
391 genetic evaluation of Holstein final score. *J. Dairy Sci.* 93: 743–752.
- 392 Aguilar, I., I. Misztal, A. Legarra, and S. Tsuruta. 2011. Efficient computation of the genomic
393 relationship matrix and other matrices used in single-step evaluation. *J. Anim. Breed.*
394 *Genet.* 128: 422–428.
- 395 Aguilar, I., A. Legarra, S. Tsuruta, and I. Misztal. 2013. Genetic evaluation using unsymmetric
396 single step genomic methodology with large number of genotypes *Interbull Bull.* 47:
397 222–225.
- 398 Bijma, P. 2012. Accuracies of estimated breeding values from ordinary genetic evaluations do
399 not reflect the correlation between true and estimated breeding values in selected
400 populations. *J. Anim. Breed. Genet.* 129: 345–358.
- 401 Christensen, O. F., and M. S. Lund. 2010. Genomic prediction when some animals are not
402 genotyped. *Genet. Sel. Evol.* 42: 2.
- 403 Fernando, R. L., J. C. M. Dekkers, and D. Garrick. 2014. A class of Bayesian methods to
404 combine large numbers of genotyped and non-genotyped animals for whole-genome
405 analyses. *Genet. Sel. Evol.* 46: 50.
- 406 Fragomeni, B. O., D. A. L. Lourenco, S. Tsuruta, Y. Masuda, I. Aguilar, A. Legarra, T. J.
407 Lawlor, and I. Misztal. 2015. Use of genomic single-step genomic BLUP with a large
408 number of genotypes. *J. Dairy Sci.* (*accepted*).
- 409 Garrick, D. J., J. F. Taylor, and R. L. Fernando. 2009. Deregressing estimated breeding values
410 and weighting information for genomic regression analyses. *Genet. Sel. Evol.* 41: 55.

- 411 Henderson, C. R. 1975. Best linear unbiased estimation and prediction under a selection model.
412 Biometrics 31: 423–447.
- 413 Legarra, A., C. R. Granie, E. Manfredi, and J. M. Elsen. 2008. Performance of genomic selection
414 in mice. Genetics 180: 611–618.
- 415 Legarra, A., and V. Ducrocq. 2012. Computational strategies for national integration of
416 phenotypes, genomic, and pedigree data in single-step best linear unbiased prediction. J.
417 Dairy Sci. 95: 4629–4645.
- 418 Legarra, A., O. F. Chistensen, I. Aguilar, and I. Misztal. 2014. Single step, a general approach
419 for genomic selection. Livest. Prod. Sci. 166: 54–65.
- 420 Liu, Z., M. E. Goddard, F. Reinhardt, and R. Reents. 2014. A single-step genomic model with
421 direct estimation of marker effect. J. Dairy Sci. 97: 5833–5850.
- 422 Lourenco, D. A. L., I. Misztal, H. Wang, I. Aguilar, S. Tsuruta, and J. K. Bertrand. 2013.
423 Prediction accuracy for a simulated maternally affected trait of beef cattle using different
424 genomic evaluation models. J. Anim. Sci. 91: 4090–4098.
- 425 Lourenco, D. A. L., I. Misztal, B. O. Fragomeni, S. Tsuruta, I. Aguilar, B. Zumbach, R. J.
426 Hawken, and A. Legarra. 2015. Accuracies of males and females with genomic
427 information on males, females, or both: a broiler chicken example. Genet. Sel. Evol.
428 *(Submitted)*.
- 429 Meuwissen, T. H. E., B. J. Hayes, and M. E. Goddard. 2001. Prediction of total genetic value
430 using genome-wide dense marker maps. Genetics 157: 1819–1829.
- 431 Misztal, I., A. Legarra, and I. Aguilar. 2014. Using recursion to compute the inverse of the
432 genomic relationship matrix. J. Dairy Sci. 97: 3943–3952.

- 433 Patry, C., and V. Ducrocq. 2011. Evidence of biases in genetic evaluations due to genomic
434 preselection in dairy cattle. *J. Dairy Sci.* 94: 1011–1020.
- 435 Pintus, M. A., E. L. Nicolazzi, J. B. C. H. M. V. Kaam, S. Biffani, A. Stella, G. Gaspa, C.
436 Dimauro, and N. P. P. Macciotta. 2013. Use of different statistical models to predict
437 direct genomic values for productive and functional traits in Italian Holsteins *J. Anim.*
438 *Breed. Genet.* 130: 32–40.
- 439 Ramirez-Valverde, R., I. Misztal, and J. K. Bertrand. 2001. Comparison of threshold vs linear
440 and animal vs sire models for predicting direct and maternal effects on calving difficulty
441 in beef cattle. *J. Anim. Sci.* 79: 333–338.
- 442 Su, G., O. F. Christensen, L. Janss, and M. S. Lund. 2014. Comparison of genomic predictions
443 using genomic relationship matrices built with different weighting factors to account for
444 locus-specific variances. *J. Dairy Sci.* 97: 6547–6559.
- 445 Sun, X., R. L. Fernando, D. J. Garrick, and J. C. M. Dekkers. 2011. An iterative approach for
446 efficient calculation of breeding values and genome-wide association analysis using
447 weighted genomic BLUP. *J. Anim. Sci.* 89 (E-Suppl.2): 28. (Abstr.).
- 448 VanRaden, P. M., and G. R. Wiggans. 1991. Deviation, calculation, and use of national animal
449 model information. *J. Dairy Sci.* 74: 2737–2746.
- 450 VanRaden, P. M. 2008. Efficient methods to compute genomic predictions. *J. Dairy Sci.* 91:
451 4414–4423.
- 452 VanRaden, P. M., C. P. VanTassel, G. R. Wiggans, T. S. Sonstegard, R. D. Schnabel, J. F.
453 Taylor, and F. S. Schenkel. 2009. Invited review: Reliability of genomic predictions for
454 North American Holstein bulls. *J. Dairy Sci.* 92: 16–24.

- 455 VanRaden, P. M., J. R. O'Connell, G. R. Wiggans, and K. A. Weigel. 2011. Genomic evaluations
456 with many more genotypes. *Genet. Sel. Evol.* 43: 10.
- 457 VanRaden, P. M., J. R. Wright, and T. A. Cooper. 2012. Adjustment of selection index
458 coefficients and polygenic variance to improve regressions and reliability of genomic
459 evaluations. *J. Dairy Sci.* 95 (Suppl. 2): 520 (Abstr.).
- 460 VanRaden, P. M., and J. R. Wright. 2013. Measuring genomic pre-selection in theory and in
461 practice. *Interbull Bull.* 47: 147–150.
- 462 Wang, H., I. Misztal, I. Aguilar, A. Legarra, and W. M. Muir. 2012. Genome-wide association
463 mapping including phenotypes from relatives without genotypes. *Genet. Res.* 94: 73–83.
- 464 Wang, H., I. Misztal, I. Aguilar, A. Legarra, R. L. Fernando, Z. Vitezica, R. Okimoto, T. Wing,
465 R. Hawken, and W. M. Muir. 2014. Genome-wide association mapping including
466 phenotypes from relatives without genotypes in a single-step (ssGWAS) for 6-week body
467 weight in broiler chickens. *Front. Genet.* 5: 134.
- 468 Wiggans, G. R., P. M. VanRaden, and T. A. Cooper. 2014. Technical note: Rapid calculation of
469 genomic evaluations for new animals. *J. Dairy Sci.* 98: 2039–2042.
- 470 Winkelman, A. M., D. L. Johnson, and B. L. Harris. 2015. Application of genomic evaluation to
471 dairy cattle in New Zealand. *J. Dairy Sci.* (*In press*).

472

473 **Table1.** Heritability (h^2) and general statistics for growth traits and CE

Trait ¹	h^2	Number of records	Average (kg)	SD (kg)	Number of genotyped animals with records
BW	0.41	6,189,661	36.47	4.45	50,784
WW	0.20	6,890,625	263.13	44.63	51,830
PWG	0.20	3,387,252	162.25	67.00	36,196
CE	0.12	1,310,684	-	-	10,558
easy	-	1,215,571	-	-	10,228
difficult	-	95,113	-	-	330

474 ¹ BW = birth weight; WW = weaning weight; PWG = post-weaning gain; CE = calving ease.

475 **Table 2.** Predictive ability of future phenotypes for young genotyped animals born in 2013

Trait ¹	Animals in validation	BLUP	ssGBLUP ²		
			ref_bulls ³	ref_2k	ref_33k
BW	18,721	0.29	0.34	0.34	0.39
WW	18,721	0.34	0.35	0.35	0.38
PWG	18,721	0.23	0.27	0.27	0.29
CE	13,166	0.12	0.13	0.13	0.13

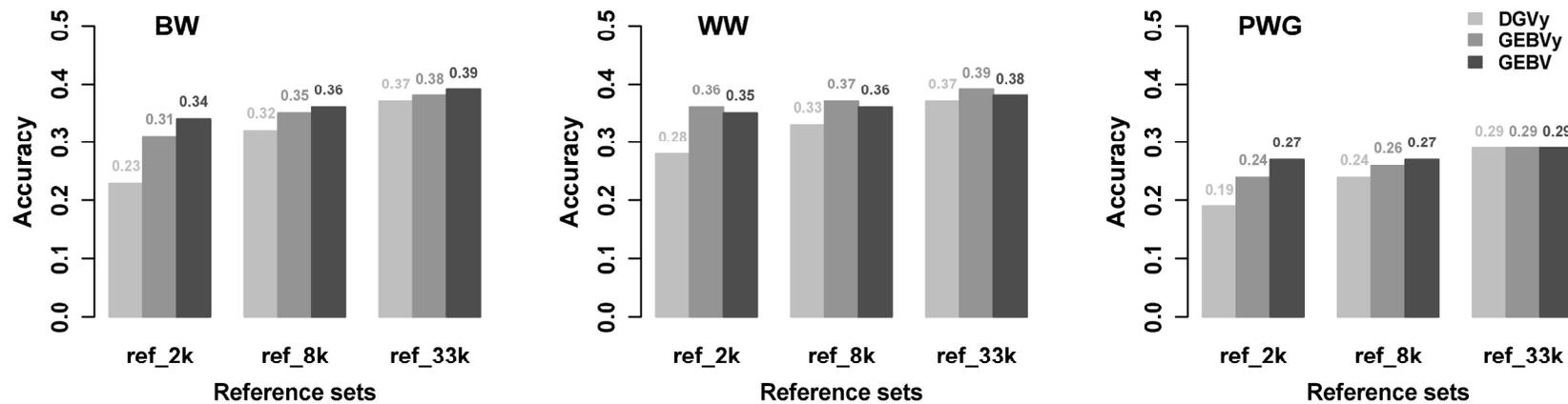
476 ¹ BW = birth weight; WW = weaning weight; PWG = post-weaning gain; CE = calving ease.477 ² Single-step genomic BLUP (ssGBLUP) included genotypes for reference and validation
478 populations, but phenotypes for validation animals were removed. Predictive ability was
479 calculated as correlation between corrected phenotypes and genomic EBV.480 ³ ref_bulls is a reference populations that contains top bulls, ref_2k contains top bulls and top
481 cows, and ref_33k contains all genotyped animals born up to 2012.

482 **Table 3.** Correlations between GEBV from full ssGBLUP and DGV_y or $GEBV_y$ from indirect
 483 predictions.

Trait	Indirect Prediction ¹	ref_2k ²	ref_8k	ref_33k
BW	DGV_y	0.66	0.87	0.96
	$GEBV_y$	0.85	0.94	0.97
WW	DGV_y	0.75	0.89	0.95
	$GEBV_y$	0.90	0.95	0.97
PWG	DGV_y	0.78	0.90	0.96
	$GEBV_y$	0.91	0.96	0.97

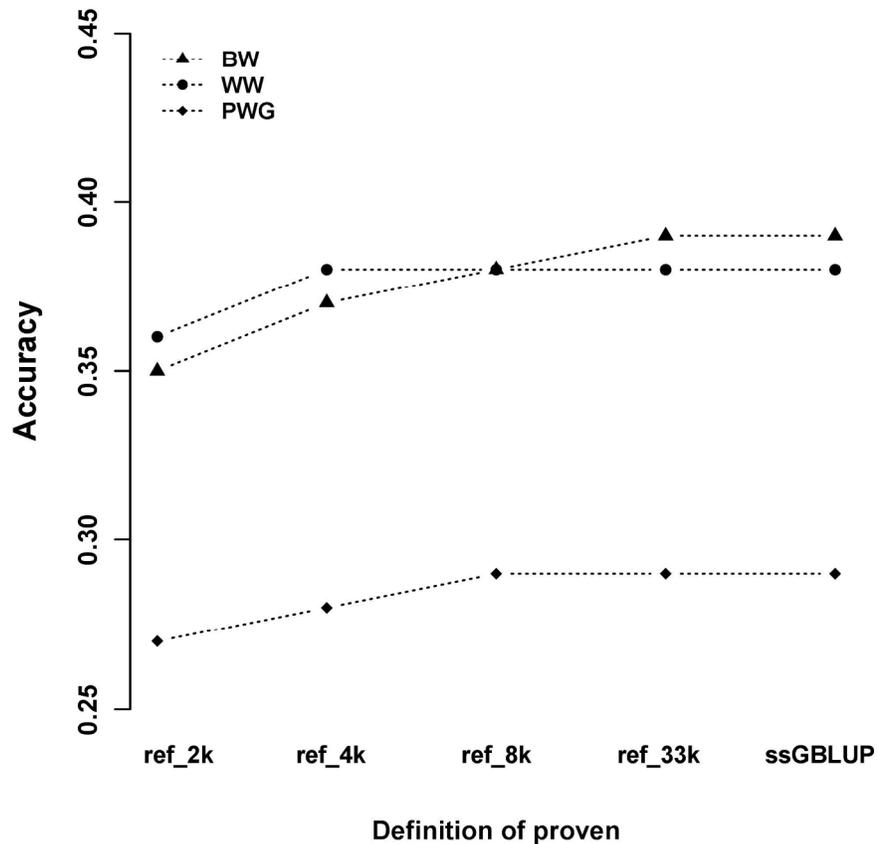
484 ¹ DGV_y is direct genomic value; $GEBV_y$ is the indirect genomic EBV obtained by an index
 485 combining parent average and DGV_y .

486 ² ref_2k is a reference populations that contains top bulls and top cows, ref_8k contains all
 487 parents that were genotyped, and ref_33k contains all genotyped animals born up to 2012.



488 **Figure 1.** Predictive ability of indirect predictions on 18,721 young genotyped animals when using reference populations ref_2k,
 489 ref_8k, and ref_33k animals to run single-step genomic BLUP (ssGBLUP) and derive SNP effects; ref_2k is a reference populations
 490 that contains top bulls and top cows, ref_8k contains all parents that were genotyped, and ref_33k contains all genotyped animals born
 491 up to 2012. DGV_y is direct genomic value; GEBV_y is the indirect genomic EBV obtained by an index combining parent average and
 492 DGV_y; GEBV is genomic predictions obtained directly from ssGBLUP when genotypes on reference and validation animals were
 493 considered together in evaluations. Predictive ability was calculated as correlation between corrected phenotypes and genomic EBV.

494



495

496 **Figure 2.** Predictive ability of GEBV for 18,721 young genotyped animals when using APY
 497 (algorithm for proven and young animals) to invert **G** matrix (genomic-based relationship
 498 matrix) with different definitions of proven animals: ref_2k, ref_4k, ref_8k, and ref_33k; ref_2k
 499 is a reference populations that contains top bulls and top cows, ref_4k contains genotyped
 500 parents of genotyped animals, ref_8k contains all parents that were genotyped, and ref_33k
 501 contains all genotyped animals born up to 2012. Predictive ability was calculated as correlation
 502 between corrected phenotypes and genomic EBV. Predictions in single-step genomic BLUP
 503 (ssGBLUP) are obtained through direct inversion of **G**.