

# Impact of genomic preselection on subsequent genetic evaluations with ssGBLUP - using real data from pigs

Ibrahim Jibrila (✉ [ibrahim.jibrila@wur.nl](mailto:ibrahim.jibrila@wur.nl))

Wageningen University & Research <https://orcid.org/0000-0002-5683-1263>

**Jeremie Vandenplas**

Wageningen University & Research

**Jan ten Napel**

Wageningen University & Research

**Rob Bergsma**

Topigs Norsvin Research Center B. V.

**Roel F Veerkamp**

Wageningen University & Research

**Mario P.L Calus**

Wageningen University & Research

---

## Research Article

**Keywords:** genomic preselection, genetic evaluation, single-step genomic best linear unbiased prediction (ssGBLUP)

**Posted Date:** June 21st, 2021

**DOI:** <https://doi.org/10.21203/rs.3.rs-638665/v1>

**License:** © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

1 **Impact of genomic preselection on subsequent genetic**  
2 **evaluations with ssGBLUP - using real data from pigs**

3 *Ibrahim Jibrila<sup>1\*</sup>, Jeremie Vandenplas<sup>1</sup>, Jan ten Napel<sup>1</sup>, Rob Bergsma<sup>2</sup>, Roel F Veerkamp<sup>1</sup> and*  
4 *Mario P.L Calus<sup>1</sup>*

5 <sup>1</sup> Wageningen University and Research Animal Breeding and Genomics, PO Box 338  
6 6700 AH Wageningen, the Netherlands

7 <sup>2</sup> Topigs Norsvin Research Center B.V., Schoenaker 6, 6641 SZ Beuningen, the  
8 Netherlands

9 \*Corresponding author

10 Email addresses and ORCID:

11 IJ: [ibrahim.jibrila@wur.nl](mailto:ibrahim.jibrila@wur.nl); ORCID: 0000-0002-5683-1263

12 JV: [jeremie.vandenplas@wur.nl](mailto:jeremie.vandenplas@wur.nl); ORCID: 0000-0002-2554-072X

13 JtN: [jan.tennapel@wur.nl](mailto:jan.tennapel@wur.nl); ORCID: 0000-00002-1918-9080

14 RB: [rob.bergsma@topignorsvin.com](mailto:rob.bergsma@topignorsvin.com); ORCID: 0000-0002-8254-5535

15 RFV: [roel.veerkamp@wur.nl](mailto:roel.veerkamp@wur.nl); ORCID: 0000-0002-5240-6534

16 MPLC: [mario.calus@wur.nl](mailto:mario.calus@wur.nl); ORCID: 0000-0002-3213-704X

17

18

19

20

21

22

23

24

## 25 **Abstract**

### 26 **Background**

27 Empirically assessing the impact of preselection on subsequent genetic evaluations of  
28 preselected animals requires comparison of scenarios with and without preselection. However,  
29 preselection almost always takes place in animal breeding programs, so it is difficult, if not  
30 impossible, to have a dataset without preselection. Hence most studies on preselection used  
31 simulated datasets, concluding that subsequent genomic estimated breeding values (GEBV)  
32 from single-step genomic best linear unbiased prediction (ssGBLUP) are unbiased. The aim of  
33 this study was to investigate the impact of genomic preselection, using real data, on accuracy  
34 and bias of GEBV of validation animals.

### 35 **Methods**

36 We used data on four pig production traits from one sire-line and one dam-line, with more  
37 intense original preselection in the dam-line than in the sire-line. The traits are average daily  
38 gain during performance testing, average daily gain throughout life, backfat, and loin depth. Per  
39 line, we ran ssGBLUP with the entire data until validation generation and considered this  
40 scenario as the reference scenario. We then implemented two scenarios with additional layers  
41 of genomic preselection by removing all animals without progeny either i) only in the validation  
42 generation, or ii) in all generations. In computing accuracy and bias, we compared GEBV  
43 against progeny yield deviation of validation animals.

### 44 **Results**

45 Results showed only a limited loss in accuracy due to the additional layers of genomic  
46 preselection. This is true in both lines, for all traits, and regardless of whether validation animals  
47 had records or not. Bias too was largely absent, and did not differ greatly among corresponding  
48 scenarios with or without additional layers of genomic preselection.

### 49 **Conclusion**

50 We concluded that impact of recent and/or historical genomic preselection is minimal on  
51 subsequent genetic evaluations of selection candidates, if these subsequent genetic evaluations  
52 are done using ssGBLUP.

53

54

55

56

57

## 58 **Background**

59 In animal breeding, parents of the next generation are often selected in multiple stages, and the  
60 initial stages of this selection are called preselection [1–3]. Selection candidates that survive  
61 preselection are called preselected animals [1–3], and those that do not are called preculled  
62 animals [3,4]. Preselection aims to reduce costs and efforts spent on animals that are not  
63 interesting for the breeding program, and achieves this by avoiding phenotyping or further  
64 testing of the preculled animals. As preculled animals have neither progeny nor records for  
65 some or all breeding goal traits, they are generally not included in subsequent genetic  
66 evaluations (i.e. genetic evaluations that come after preselection). Preselection therefore  
67 decreases the amount of information available for subsequent genetic evaluations of preselected  
68 animals. Properly assessing the impact of preselection on subsequent genetic evaluation of  
69 preselected animals requires a scenario without preselection, against which scenarios with  
70 preselection can be compared. Because in animal breeding programmes preselection almost  
71 always takes place, it is difficult, if not impossible, to have a scenario without preselection. This  
72 is why most studies available on preselection used simulated datasets [e.g. 1–5]. Those studies  
73 have shown that when a subsequent genetic evaluation of preselected animals is done using  
74 pedigree-based best linear unbiased prediction (PBLUP), genomic preselection results in  
75 accuracy loss and bias in the estimated breeding values (EBV) of preselected animals [1,6–9].  
76 Some of these studies [6–9] further showed that the accuracy loss and bias caused by genomic  
77 preselection can be avoided if the information on preculled animals that was utilized at  
78 preselection is included in the subsequent PBLUP evaluation. On the other hand, our previous  
79 works [3,4] have shown that when the subsequent genetic evaluation is done with single-step  
80 genomic BLUP (ssGBLUP), genomic EBV (GEBV) of preselected animals are estimated  
81 without bias. We [4] further showed that to avoid genomic preselection bias in subsequent

82 ssGBLUP evaluation of preselected animals, genotypes of their preculled sibs are only needed  
83 if not all of their parents are genotyped.

84 In our previous works [3,4], being based on simulated datasets, preselection was the only  
85 possible source of bias in ssGBLUP evaluations. However, in real breeding programmes, other  
86 sources of bias in ssGBLUP evaluations may exist and are potentially difficult to control.  
87 Therefore, impact of preselection might be confounded by the impact of these other factors.  
88 These other possible sources of bias include, amongst others, inaccurate or incomplete pedigree  
89 [10], inaccurately estimated additive genetic (co)variances [10], and a reference population of  
90 selected genotyped animals [11,12]. Although some ways of reducing the bias caused by these  
91 factors have been developed, the bias is usually not completely eliminated in evaluations using  
92 real data (e.g. [10–12]). This may explain the observation that in practice GEBV obtained from  
93 ssGBLUP evaluations are sometimes biased. The aim of this study was to investigate the impact  
94 of genomic preselection on subsequent ssGBLUP evaluations, using real data from an ongoing  
95 pig breeding program in which preselection has taken place. To achieve this aim, we used the  
96 full dataset as control and retrospectively implemented additional layers of genomic  
97 preselection, and results from subsequent ssGBLUP evaluations after these additional layer of  
98 genomic preselection were compared against results from ssGBLUP evaluation of the full  
99 available data.

## 100 **Methods**

### 101 **Data**

102 In our analyses, additional layers of genomic preselection were implemented when the animals  
103 already had phenotypes, by discarding animals that did not have progeny in the data. Our  
104 subsequent genetic evaluations only involved reevaluating preselected animals, either with or  
105 without preculled animals in the reevaluations. We separated the available data in two parts,

106 according to a cut-off birth date. Animals born before or on the cut-off birth date were used as  
107 reference population, and animals born after the cut-off birth date were used as validation  
108 population, from which animals were selected to be used for validation (these are hereafter  
109 referred to as “validation animals”). Only animals in the validation population that met the  
110 following two requirements were selected as validation animals: 1) none of their parents were  
111 included in the validation population, and 2) they had progeny associated with phenotypes. The  
112 first requirement ensured that validation animals represented the youngest generation of  
113 selection candidates in a breeding program in practice, and not multiple generations. The second  
114 requirement enabled validation of the GEBV of the validation animals against their progeny  
115 yield deviation (PYD) [13]. Meeting the second requirement was needed, because own  
116 phenotypes of the validation animals were used in our subsequent evaluations, and could thus  
117 not be used to validate their GEBV.

118 We obtained pig production traits data on one sire-line and one dam-line from Topigs Norsvin.  
119 These data were collected between 1970 and 2020, and the traits were average daily gain during  
120 performance testing, average daily gain throughout the lifetime, backfat, and loin depth. Topigs  
121 Norsvin (pre)selected both lines on these production traits. However, there was more emphasis  
122 on reproduction traits than on production traits in the dam-line. Details on the amount of data  
123 utilized in this study are in Table 1. The data were recorded on originally preselected animals  
124 (i.e. the animals preselected by Topigs Norsvin), with the sire-line being much more balanced  
125 than the dam-line, in terms of proportions males and females with records per generation (ratio  
126 of males with records to females with records is about 50:50 in the sire-line and about 20:80 in  
127 the dam-line). We studied impact of genomic preselection in the two lines separately, because  
128 the traits we studied had different weights in breeding goals of the two lines. The cut-off date  
129 to split the data into reference and validation populations was 31<sup>st</sup> January, 2017 for the sire

130 line, and 31<sup>st</sup> December, 2015 for the dam-line. In the pedigree, animals with one or both parents  
131 missing were assigned to genetic groups, according to line and year of birth of each animal.

### 132 **Genomic data and quality control**

133 Our genomic data included genotypes of animals for about 21,000 SNP segregating in both  
134 lines, and distributed across the 18 autosomes in the pig genome. The SNP were genotyped  
135 using a custom SNP chip. We used Plink [14] for all quality control operations on our genomic  
136 data. Per genomic preselection scenario (as described later) and per line, animals and SNPs with  
137 call rates less than 90% were removed, as well as SNPs that deviated from Hardy-Weinberg  
138 equilibrium (Hardy-Weinberg equilibrium exact test p value =  $10^{-15}$ ), or had a minor allele  
139 frequency below 0.005. Table 1 contains the summary of the pedigree, genomic and phenotypic  
140 information utilized in the subsequent genetic evaluations following each genomic preselection  
141 scenario, per line.

### 142 **Computation of pre-corrected phenotypes**

143 In our genetic evaluations, we used pre-corrected phenotypes (rather than raw phenotypes) as  
144 records. Animals of different lines were sometimes raised together, so they shared some fixed  
145 and non-genetic random effects. Because we studied impact of genomic preselection within  
146 lines, it was necessary to correct phenotypes for all non-genetic effects before the data was  
147 divided into lines. Another motivation for using pre-corrected phenotypes was that some classes  
148 of these non-genetic effects could include only one or a few animals per class due to our  
149 implemented additional preselection. We used the following multi-trait pedigree-based animal  
150 model to compute pre-corrected phenotypes for all traits:

$$151 \quad \mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{W}\mathbf{p} + \mathbf{Z}\mathbf{u} + \mathbf{e}, \quad (\text{eq. 1}),$$

152 where  $\mathbf{y}$  was the vector of phenotypes;  $\mathbf{b}$  was the vector of fixed effects, with incidence matrix  
153  $\mathbf{X}$ ;  $\mathbf{p}$  was the vector of non-genetic random effects, with incidence matrix  $\mathbf{W}$ ;  $\mathbf{u}$  was the vector

154 of breeding values, with incidence matrix  $\mathbf{Z}$ ; and  $\mathbf{e}$  was the vector of residuals. Then for every  
155 animal (i) with phenotype, precorrected phenotype ( $y_{ci}$ ) was:

$$156 \quad y_{ci} = \hat{u}_i + \hat{e}_i \quad (\text{eq. 2}).$$

157 The (co)variance components used for this analysis were estimated, before separating the data  
158 into lines, from a multi-trait pedigree-based animal model in ASReml [15] using **eq. 1**. All  
159 computations of (G)EBV were performed using MiXBLUP [16].

### 160 **Preselection**

161 Per line, we implemented a reference scenario and two scenarios that added layers of genomic  
162 preselection. The reference scenario - against which other scenarios could be compared - only  
163 included the original genomic preselection implemented by Topigs Norsvin. Thus, the  
164 subsequent ssGBLUP evaluations following the reference scenario utilized the entire available  
165 data until the validation generation. The second scenario is called validation generation  
166 preselection (the VGP scenario). In this scenario, we only implemented additional genomic  
167 preselection in the validation generation, by discarding all animals in the validation generation  
168 that had no progeny in the data, but had genotypes and/or phenotypes. This scenario was  
169 implemented to study the impact of extreme genomic preselection in a single generation. The  
170 third scenario is called multi-generation preselection (the MGP scenario), in which we  
171 discarded any animal in the validation and previous generations with no progeny in the data.  
172 This scenario was implemented to study the carry-over impact of extreme genomic preselection  
173 in multiple generations. Animals kept after each of the genomic preselection scenarios are  
174 shown in Figure 1.

### 175 **Subsequent genetic evaluations**

176 Following every scenario of genomic preselection, we implemented a subsequent ssGBLUP

177 evaluation with all animals that survived the genomic preselection. We call this evaluation  
178 subsequent because it came after the initial evaluation that provided the GEBV used in  
179 preselection. The ssGBLUP evaluations were conducted using MiXBUP [16], with and  
180 without records (i.e. own phenotypes) on the animals in the validation generation (see Table 1).  
181 Progeny of validation animals were not included in the subsequent genetic evaluations. We  
182 estimated variance components after every preselection scenario, per line, using a pedigree-  
183 based multi-trait animal model in ASReml. We used these scenario-specific variance  
184 components in the subsequent genetic evaluations to ensure that the variance components used  
185 were appropriate for the pre-corrected phenotypes. At the subsequent genetic evaluations, the  
186 model used for the estimations of both variance components and breeding values was:

$$187 \quad \mathbf{y} = \mathbf{x}\mathbf{b} + \mathbf{Z}\mathbf{u} + \mathbf{e} \quad (\text{eq. 3}),$$

188 where  $\mathbf{y}$  was the vector of pre-corrected phenotypes;  $\mathbf{x}$  and  $\mathbf{Z}$  were incidence vector and matrix  
189 linking pre-corrected phenotypes to overall mean and random animal effects, respectively;  $\mathbf{b}$   
190 was the overall mean;  $\mathbf{u}$  was the vector of breeding values; and  $\mathbf{e}$  was the vector of residuals.  
191 We also repeated all subsequent genetic evaluations using PBLUP, to verify the impact of using  
192 genotypes on the observed results.

193 **Figure 1 here**

194  
195  
196  
197  
198  
199  
200  
201  
202  
203  
204  
205  
206  
207

208 **Table 1** Data utilized in subsequent ssGBLUP<sup>a</sup> evaluations following each preselection  
 209 scenario, after quality control

Data in the subsequent ssGBLUP evaluation/Preselection scenario	With records on animals in the validation generation			Without records on animals in the validation generation		
	Reference <sup>b</sup>	VGP <sup>c</sup>	MGP <sup>d</sup>	Reference <sup>b</sup>	VGP <sup>c</sup>	MGP <sup>d</sup>
<i>The sire line</i>						
Number of animals in the pedigree	81,875	60,950	12,777	81,875	60,950	12,777
Number of animals with record for at least one trait	75,129	54,217	6,065	52,846	52,846	4,694
Number of animals with genotypes	33,506	23,315	5,131	33,506	23,315	5,131
Number of SNP	20,550	20,963	20,926	20,550	20,963	20,926
<i>The dam line</i>						
Number of animals in the pedigree	160,426	124,031	33,485	160,426	124,031	33,485
Number of animals with record for at least one trait	139,403	103,018	12,514	100,710	100,710	10,206
Number of animals with genotypes	50,895	36,369	9,072	50,895	36,369	9,072
Number of SNP	19,199	19,256	20,647	19,199	19,256	20,647

210 <sup>a</sup> single-step genomic best linear unbiased prediction

211 <sup>b</sup> In the reference scenario, the subsequent ssGBLUP evaluation utilized the entire available  
 212 data until the validation generation

213 <sup>c</sup> Validation generation preselection (VGP) scenario. In this scenario, additional genomic  
 214 preselection was only implemented in the validation generation, by discarding all animals in  
 215 the validation generation that did not have progeny in the data.

216 <sup>d</sup> Multi-generation preselection (MGP) scenario. In this scenario, any animal in the validation  
 217 or reference generations with no progeny in the data was discarded.

## 218 Implementation of single-step GBLUP

219 The inverse of the combined pedigree-genomic relationship ( $\mathbf{H}^{-1}$ ) was obtained as follows

220 [17,18]:

$$221 \quad \mathbf{H}^{-1} = \mathbf{A}^{-1} + \begin{bmatrix} 0 & 0 \\ 0 & (0.95\mathbf{G}_t + 0.05\mathbf{A}_{22})^{-1} - \mathbf{A}_{22}^{-1} \end{bmatrix} \quad (\text{eq. 4}),$$

222 where  $\mathbf{A}^{-1}$  was the inverse of the pedigree relationship matrix, and  $\mathbf{A}_{22}$  was part of the pedigree

223 relationship matrix referring to genotyped animals. We considered inbreeding in setting up both

224  $\mathbf{A}^{-1}$  and  $\mathbf{A}_{22}$  to avoid bias caused by ignoring inbreeding (Tsuruta et al., 2019). The genomic

225 relationship matrix  $\mathbf{G}_t$  was computed as follows:

$$226 \quad \mathbf{G}_t = (1 - \bar{f}_p)\mathbf{G}_r + 2\bar{f}_p\mathbf{1}\mathbf{1}' \quad (\text{eq. 5}),$$

227 where  $\bar{f}_p$  was the average pedigree inbreeding coefficient across genotyped animals,  $\mathbf{G}_r$  was the  
228 raw genomic relationship matrix computed following the first method of VanRaden [19], and  
229  $\mathbf{1}\mathbf{1}'$  was a matrix of 1s. The scaling of  $\mathbf{G}_r$  to  $\mathbf{G}_t$  was done to make the average genomic  
230 inbreeding equal to the average pedigree inbreeding, i.e. to have  $\mathbf{G}$  and  $\mathbf{A}_{22}$  on the same scale  
231 so that they are compatible. As the animals with genotypes in this study were selectively  
232 genotyped, this transformation made sure that the impact of selective genotyping was taken  
233 care of [11,12]. In computing  $\mathbf{G}_r$ , we computed (current) allele frequencies using all available  
234 genomic data after quality control. We gave the weights of 0.95 to  $\mathbf{G}_t$  and 0.05 to  $\mathbf{A}_{22}$  to ensure  
235 that  $\mathbf{G}$  was invertible [17,18].

### 236 **Measures of accuracy and bias in the subsequent genetic evaluations**

237 We used progeny yield deviation (PYD) [13] as a proxy for true breeding value (TBV), against  
238 which GEBV were compared when computing accuracy and bias. To compute PYD, we ran a  
239 multi-trait pedigree-based animal model per line in MiXBLUP, with precorrected phenotypes  
240 as records and an overall mean as the only fixed effect (eq. 3). The (co)variance components  
241 used in this model were also estimated per line in ASReml, from precorrected phenotypes using  
242 a multi-trait pedigree-based animal model that only included a mean fixed effect (eq. 3). From  
243 the output of this analysis, we computed PYD for each trait for all validation sires and dams as:

$$244 \quad PYD_i = \frac{\sum_{p=1}^n y_{cp} - g_m}{n} \quad (\text{eq. 6}),$$

245 where  $PYD_i$  was the progeny yield deviation of a sire or dam  $i$ ,  $y_{cp}$  was the precorrected  
246 phenotype of a progeny  $p$  of the sire or dam  $i$ ,  $g_m$  was the genetic contribution of the mate of

247 sire or dam  $i$  to  $y_{cp}$ , and  $n$  was the number of phenotyped progeny of sire or dam  $i$ . Estimation  
 248 of PYD was done before discarding progeny of validation animals from the data. Since progeny  
 249 of validation animals were not included in subsequent genetic evaluations, comparing (G)EBV  
 250 to PYD can be considered as a forward-in-time validation. To account for differences in number  
 251 of progeny used in estimating PYD for different validation animals when estimating accuracy  
 252 and bias, we approximated the reliability of PYD for each validation animal for each trait as:

$$253 \quad \frac{1/4nh^2}{1+1/4(n-1)h^2} \quad (\text{eq. 7}),$$

254 where  $n$  was the validation animal's number of half-sib progeny with records, and  $h^2$  was the  
 255 heritability of the trait [20]. For convenience, we assumed all progeny of a validation animal  
 256 were half-sibs, though some of them were full-sibs.

257 Validation accuracy was computed as weighted Pearson's correlation coefficient between  
 258 PYD and GEBV of all validation animals, with reliability of PYD used as the weight. We  
 259 computed two types of bias. The first type is absolute bias, which is a measure of whether  
 260 estimated genetic gain is equal to true genetic gain. Absolute bias was computed as the  
 261 weighted mean difference between PYD and half of the (G)EBV of all validation animals,  
 262 expressed in additive genetic standard deviation (SD) units of the trait. A negative difference  
 263 means that GEBV are on average overestimated, and therefore genetic gain is overestimated,  
 264 and vice versa. Before computing differences between PYD and half of the (G)EBV of  
 265 validation animals, we made sure that PYD and (G)EBV were on the same scale. We did this  
 266 in the following steps: from the model used in computing PYD, we computed average EBV  
 267 across all animals in the first three reference generations. We then subtracted half of this  
 268 average EBV from PYD of each validation animal. Then from each subsequent genetic  
 269 evaluation, we computed the average (G)EBV of all animals in the first three reference

270 generations. We then subtracted this average (G)EBV from (G)EBV of each validation animal.  
271 The second type of bias we computed is dispersion bias. Dispersion bias was measured by the  
272 weighted regression coefficient of PYD on (G)EBV of all validation animals. If the regression  
273 coefficient is equal to the expected value, then there is no dispersion bias. Note that the  
274 expected value is 0.5, because PYD only includes half of the breeding value of a parent. A  
275 regression coefficient less than the expected value means that variance of (G)EBV is inflated,  
276 and vice versa.

## 277 **Results**

278 Results of the subsequent genetic evaluations conducted with ssGBLUP are presented in Tables  
279 2 and 3 for the sire-line and the dam-line, respectively. Results in Tables 4 and 5 are from  
280 subsequent genetic evaluations done with PBLUP, respectively for the sire-line and the dam-  
281 line. In addition to validation accuracy and bias, we also showed the estimated heritability for  
282 every subsequent genetic evaluation scenario, and number of validation animals.

283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304

305 **Table 2** Performance of ssGBLUP<sup>a</sup> in the subsequent genetic evaluations in the sire-line

Measure/Preselection scenario	With records on animals in the validation generation			Without records on animals in the validation generation		
	Reference <sup>b</sup>	VGP <sup>c</sup>	MGP <sup>d</sup>	Reference <sup>b</sup>	VGP <sup>c</sup>	MGP <sup>d</sup>
<i>Average daily gain during performance testing, number of validation animals = 1382</i>						
Estimated heritability	0.24	0.25	0.33	0.24	0.24	0.35
Validation accuracy	0.51	0.51	0.50	0.47	0.47	0.44
Absolute bias	-0.09	-0.15	-0.01	-0.11	-0.11	-0.02
Dispersion bias	0.48	0.49	0.48	0.48	0.48	0.46
<i>Average daily gain throughout life, number of validation animals = 1383</i>						
Estimated heritability	0.26	0.28	0.33	0.27	0.27	0.35
Validation accuracy	0.57	0.56	0.55	0.52	0.52	0.48
Absolute bias	-0.10	-0.17	-0.06	-0.14	-0.14	-0.08
Dispersion bias	0.48	0.49	0.50	0.47	0.47	0.49
<i>Backfat, number of validation animals = 1383</i>						
Estimated heritability	0.58	0.58	0.58	0.58	0.58	0.60
Validation accuracy	0.69	0.68	0.67	0.63	0.63	0.56
Absolute bias	-0.02	-0.03	-0.03	-0.05	-0.05	-0.09
Dispersion bias	0.48	0.47	0.47	0.44	0.44	0.42
<i>Loin depth, number of validation animals = 1383</i>						
Estimated heritability	0.55	0.55	0.55	0.55	0.55	0.57
Validation accuracy	0.68	0.67	0.65	0.62	0.62	0.54
Absolute bias	0.01	0.00	0.00	0.00	0.00	-0.01
Dispersion bias	0.50	0.50	0.48	0.48	0.48	0.45

306 SEs were in the range 0.01-0.03 for estimated heritability and dispersion bias, and 0.01-0.02  
 307 for validation accuracy and absolute bias.

308 <sup>a</sup> single-step genomic best linear unbiased prediction

309 <sup>b</sup> In the reference scenario, the subsequent ssGBLUP evaluation utilized the entire available  
 310 data until the validation generation

311 <sup>c</sup> Validation generation preselection (VGP) scenario. In this scenario, additional genomic  
 312 preselection was only implemented in the validation generation, by discarding all animals in  
 313 the validation generation that did not have progeny in the data.

314 <sup>d</sup> Multi-generation preselection (MGP) scenario. In this scenario, any animal in the validation  
 315 or reference generations with no progeny in the data was discarded.

316  
 317  
 318  
 319  
 320  
 321  
 322  
 323  
 324  
 325  
 326  
 327  
 328  
 329  
 330  
 331

332 **Table 3** Performance of ssGBLUP<sup>a</sup> in the subsequent genetic evaluations in the dam-line

Measure/Preselection scenario	With records on animals in the validation generation			Without records on animals in the validation generation		
	Reference <sup>b</sup>	VGP <sup>c</sup>	MGP <sup>d</sup>	Reference <sup>b</sup>	VGP <sup>c</sup>	MGP <sup>d</sup>
<i>Average daily gain during performance testing, number of validation animals = 2323</i>						
Estimated heritability	0.31	0.32	0.40	0.30	0.30	0.38
Validation accuracy	0.35	0.31	0.29	0.28	0.28	0.23
Absolute bias	-0.05	-0.14	0.04	0.03	0.03	0.14
Dispersion bias	0.46	0.43	0.41	0.44	0.44	0.43
<i>Average daily gain throughout life, number of validation animals = 2405</i>						
Estimated heritability	0.31	0.33	0.43	0.31	0.31	0.44
Validation accuracy	0.46	0.42	0.42	0.38	0.38	0.35
Absolute bias	-0.06	-0.16	-0.01	0.00	0.00	0.08
Dispersion bias	0.45	0.42	0.42	0.43	0.43	0.43
<i>Backfat, number of validation animals = 2312</i>						
Estimated heritability	0.51	0.51	0.51	0.51	0.51	0.53
Validation accuracy	0.52	0.50	0.50	0.45	0.45	0.42
Absolute bias	0.02	-0.01	-0.03	0.02	0.02	-0.01
Dispersion bias	0.43	0.41	0.41	0.42	0.42	0.41
<i>Loin depth, number of validation animals = 1164</i>						
Estimated heritability	0.50	0.50	0.55	0.49	0.49	0.53
Validation accuracy	0.62	0.60	0.59	0.55	0.56	0.49
Absolute bias	-0.02	-0.03	0.02	-0.04	-0.04	0.03
Dispersion bias	0.54	0.54	0.52	0.53	0.53	0.51

333 SEs were in the range 0.01-0.02 for estimated heritability, validation accuracy and absolute  
 334 bias, and 0.01-0.04 for dispersion bias.

335 <sup>a</sup> single-step genomic best linear unbiased prediction

336 <sup>b</sup> In the reference scenario, the subsequent ssGBLUP evaluation utilized the entire available  
 337 data until the validation generation

338 <sup>c</sup> Validation generation preselection (VGP) scenario. In this scenario, additional genomic  
 339 preselection was only implemented in the validation generation, by discarding all animals in  
 340 the validation generation that did not have progeny in the data.

341 <sup>d</sup> Multi-generation preselection (MGP) scenario. In this scenario, any animal in the validation  
 342 or reference generations with no progeny in the data was discarded.

343

344

345

346

347

348

349

350

351

352

353 **Table 4** Performance of PBLUP<sup>a</sup> in the subsequent genetic evaluations in the sire-line

Measure/Preselection scenario	With records on animals in the validation generation			Without records on animals in the validation generation		
	Reference <sup>b</sup>	VGP <sup>c</sup>	MGP <sup>d</sup>	Reference <sup>b</sup>	VGP <sup>c</sup>	MGP <sup>d</sup>
<i>Average daily gain during performance testing, number of validation animals = 1382</i>						
Estimated heritability	0.24	0.25	0.33	0.24	0.24	0.35
Validation accuracy	0.51	0.50	0.49	0.41	0.41	0.40
Absolute bias	-0.04	-0.11	0.01	-0.01	-0.01	0.01
Dispersion bias	0.53	0.54	0.48	0.55	0.55	0.49
<i>Average daily gain throughout life, number of validation animals = 1383</i>						
Estimated heritability	0.26	0.28	0.33	0.27	0.27	0.35
Validation accuracy	0.58	0.56	0.54	0.47	0.47	0.44
Absolute bias	-0.06	-0.14	-0.04	-0.05	-0.05	-0.05
Dispersion bias	0.55	0.55	0.51	0.56	0.56	0.54
<i>Backfat, number of validation animals = 1383</i>						
Estimated heritability	0.58	0.58	0.58	0.58	0.58	0.60
Validation accuracy	0.67	0.66	0.66	0.48	0.48	0.46
Absolute bias	-0.03	-0.03	-0.03	-0.09	-0.09	-0.10
Dispersion bias	0.50	0.50	0.50	0.46	0.46	0.43
<i>Loin depth, number of validation animals = 1383</i>						
Estimated heritability	0.55	0.55	0.55	0.55	0.55	0.57
Validation accuracy	0.66	0.65	0.64	0.49	0.49	0.46
Absolute bias	0.00	0.00	0.00	0.01	0.01	0.00
Dispersion bias	0.50	0.49	0.49	0.48	0.48	0.46

354 SEs were in the range 0.01-0.03 for estimated heritability and dispersion bias, and 0.01-0.02  
 355 for validation accuracy and absolute bias.

356 <sup>a</sup> Pedigree-based best linear unbiased prediction

357 <sup>b</sup> In the reference scenario, the subsequent PBLUP evaluation utilized the entire available data  
 358 until the validation generation

359 <sup>c</sup> Validation generation preselection (VGP) scenario. In this scenario, additional genomic  
 360 preselection was only implemented in the validation generation, by discarding all animals in  
 361 the validation generation that did not have progeny in the data.

362 <sup>d</sup> Multi-generation preselection (MGP) scenario. In this scenario, any animal in the validation  
 363 or reference generations with no progeny in the data was discarded.

364

365

366

367

368

369

370

371

372

373

374

**Table 5** Performance of PBLUP<sup>a</sup> in the subsequent genetic evaluations in the dam-line

Measure/Preselection scenario	With records on animals in the validation generation			Without records on animals in the validation generation		
	Reference <sup>b</sup>	VGP <sup>c</sup>	MGP <sup>d</sup>	Reference <sup>b</sup>	VGP <sup>c</sup>	MGP <sup>d</sup>
<i>Average daily gain during performance testing, number of validation animals = 2323</i>						
Estimated heritability	0.31	0.32	0.40	0.30	0.30	0.38
Validation accuracy	0.35	0.30	0.30	0.24	0.24	0.21
Absolute bias	-0.04	-0.16	0.01	0.08	0.08	0.13
Dispersion bias	0.52	0.45	0.42	0.50	0.50	0.45
<i>Average daily gain throughout life, number of validation animals = 2405</i>						
Estimated heritability	0.31	0.33	0.43	0.31	0.31	0.44
Validation accuracy	0.48	0.43	0.43	0.34	0.34	0.31
Absolute bias	-0.05	-0.18	-0.03	0.05	0.05	0.07
Dispersion bias	0.51	0.47	0.44	0.51	0.51	0.44
<i>Backfat, number of validation animals = 2312</i>						
Estimated heritability	0.51	0.51	0.51	0.51	0.51	0.53
Validation accuracy	0.52	0.50	0.50	0.37	0.37	0.36
Absolute bias	0.02	0.00	-0.03	0.04	0.04	0.00
Dispersion bias	0.45	0.43	0.42	0.41	0.41	0.39
<i>Loin depth, number of validation animals = 1164</i>						
Estimated heritability	0.50	0.50	0.55	0.49	0.49	0.53
Validation accuracy	0.58	0.56	0.56	0.43	0.43	0.41
Absolute bias	0.00	-0.01	0.04	-0.02	-0.02	0.04
Dispersion bias	0.55	0.54	0.51	0.57	0.57	0.52

375 SEs were in the range 0.01-0.02 for estimated heritability, validation accuracy and absolute  
 376 bias, and 0.01-0.04 for dispersion bias.

377 <sup>a</sup> Pedigree-based best linear unbiased prediction

378 <sup>b</sup> In the reference scenario, the subsequent PBLUP evaluation utilized the entire available data  
 379 until the validation generation

380 <sup>c</sup> Validation generation preselection (VGP) scenario. In this scenario, additional genomic  
 381 preselection was only implemented in the validation generation, by discarding all animals in  
 382 the validation generation that did not have progeny in the data.

383 <sup>d</sup> Multi-generation preselection (MGP) scenario. In this scenario, any animal in the validation  
 384 or reference generations with no progeny in the data was discarded.

### 385 **Subsequent ssGBLUP evaluations with records on animals in the validation generation**

386 With records on animals in the validation generation included in the subsequent ssGBLUP  
 387 evaluations, estimated heritability for average daily gain traits in the sire-line increased from  
 388 the reference to validation generation preselection (VGP) to multi-generation preselection  
 389 (MGP) scenarios, with more increase from VGP to MGP than from reference to VGP. For  
 390 backfat and loin depth, the heritability remained the same across all scenarios. For the dam-

391 line, estimated heritability increased from reference to VGP to MGP scenarios, except for  
392 backfat, where it remained the same across all scenarios. Observed increases in estimated  
393 heritabilities were generally due to decreases in residual variances across the scenarios, while  
394 additive genetic variances generally remained similar (Tables S1 and S2). For both lines and  
395 for all traits, validation accuracy decreased from reference to VGP to MGP scenarios, albeit the  
396 differences were small. For both lines, absolute bias was largely absent for backfat and loin  
397 depth, and marginal for the average daily gain traits. The highest value of absolute bias recorded  
398 was -0.17 additive genetic SDs, under the VGP scenario for average daily gain throughout life  
399 in the sire-line (Table 2). Generally, the values of absolute bias for average daily gain traits  
400 moved further away from zero from reference to VGP, and then moved closest to zero with  
401 MGP. For the sire-line, regression coefficients of PYD on GEBV - an indicator of dispersion  
402 bias - showed no consistent pattern across preselection scenarios for all traits. For all traits and  
403 for all scenarios, they ranged from 0.47 to 0.50, being close to the expected value of 0.5. For  
404 the dam-line, the regression coefficients decreased or remained the same from reference to VGP  
405 to MGP scenarios. They were less than 0.5 for the two average daily gain traits and backfat.  
406 For loin depth, they were greater than 0.5.

#### 407 **Subsequent ssGBLUP evaluations without records on animals in the validation generation**

408 Without records on animals in the validation generation in the subsequent ssGBLUP  
409 evaluations, all results for the reference and VGP scenarios were the same. Just like when  
410 records on animals in the validation generation were included, here too, estimated heritability  
411 increased from reference and VGP to MGP scenarios, and in this case for all traits in both lines.  
412 Validation accuracy also decreased from reference and VGP to MGP scenarios, and in this case  
413 with bigger decreases compared to when records on animals in the validation generation were  
414 included. Absolute bias was also largely absent for backfat and loin depth for both lines, and  
415 showed no particular pattern for average daily gain traits for the two lines. Even for the average

416 daily gain traits, it was still small, with  $\pm 0.14$  additive genetic SD being the highest value  
417 (Tables 2 and 3). Regression coefficients of PYD on GEBV were similar to their corresponding  
418 value when records on animals in the validation generation were included. The only exception  
419 were all scenarios for backfat in the sire-line, where the regression coefficients of PYD on  
420 GEBV appeared to be lower than their corresponding values when records on animals in the  
421 validation generation were included. For both lines, the regression coefficients ranged from  
422 0.41 (for the MGP scenario for backfat in the dam-line) to 0.53 (for the reference and VGP  
423 scenarios for loin depth in the dam-line).

#### 424 **Subsequent genetic evaluations with PBLUP**

425 With records on animals in the validation generation included, validation accuracies from  
426 subsequent PBLUP evaluations were similar in both magnitude and pattern across the  
427 preselection scenarios and lines, to their corresponding values from subsequent ssGBLUP  
428 evaluations. However, without records on animals in the validation generation in the subsequent  
429 genetic evaluations, validation accuracies were lower with PBLUP than with ssGBLUP for all  
430 scenarios in both lines. For both lines and with or without records on animals in the validation  
431 generation, absolute bias with PBLUP was always lower than or similar to its corresponding  
432 value with ssGBLUP. Regression coefficients of PYD on (G)EBV were also bigger than or  
433 similar to their corresponding values with ssGBLUP.

## 434 **Discussion**

435 In this study, we investigated the impact of genomic preselection on subsequent ssGBLUP  
436 evaluations of preselected animals, using real data from an ongoing pig breeding program in  
437 which preselection has taken place, by retrospectively implementing additional layers of  
438 preselection. The data was on production traits of pigs from one sire-line and one dam-line. Per  
439 line, we implemented three genomic preselection scenarios. We used pre-corrected phenotypes

440 as records in the subsequent genetic evaluations, and progeny yield deviation (PYD) as the  
441 proxy for TBV. We did the subsequent genetic evaluations either with or without records on  
442 animals in the validation generation, and in all cases without progeny of validation animals. In  
443 both lines, for all traits and with or without records on validation animals, absolute bias was  
444 largely absent across the three genomic preselection scenarios, while with more preselection  
445 validation accuracy only showed small decreases and hardly any dispersion bias was induced.

446 In the two scenarios with additional genomic preselection (i.e. VGP and MGP scenarios), the  
447 preselected animals in every generation were the animals that in reality were selected and  
448 produced progeny in the next generation, and the preculled animals were those animals that  
449 were in reality culled after performance testing. Thus, these two scenarios represent either i)  
450 situations in which all the selection in a generation is done in only one stage, after selection  
451 candidates have own records, or ii) situations in which an additional selection stage is  
452 implemented after preselected animals have had progeny. While neither of these cases is true  
453 for the data we used, the scenarios we implemented enabled us to investigate the impact of  
454 genomic preselection on subsequent genetic evaluations of preselected animals using real data,  
455 by including different amounts of pedigree, genomic and phenotypic information in the  
456 subsequent genetic evaluations we implemented. The validation accuracy we computed as the  
457 correlation between (G)EBV and PYD is not numerically the same as the accuracy of predicting  
458 TBV, since variance of PYD has some non-genetic component, in addition to genetic  
459 component [13]. However, the two accuracies are proportional to each other, and this enabled  
460 us to make comparison among subsequent genetic evaluation scenarios [21].

#### 461 **Comparison of results across preselection scenarios and between ssGBLUP and PBLUP**

462 With both ssGBLUP and PBLUP, validation accuracy decreased with more genomic  
463 preselection (i.e. from reference to VGP to MGP scenarios), and this could be explained by the

464 fact that the amount of phenotypic information also reduced in that order (Table 1). In our  
465 previous study using simulated datasets [3], we found accuracy in subsequent ssGBLUP  
466 evaluations to be decreasing as amounts of phenotypic information decreased with more intense  
467 preselection. For most of the traits in the current study, estimated heritability increased with  
468 increase in genomic preselection, and this could have influenced, at least partly, the magnitude  
469 of decrease in accuracy with decrease in amount of phenotypic information due to preselection.  
470 This could also contribute to explaining why decrease in validation accuracy with more  
471 genomic preselection was small. We also observed that validation accuracy was higher with  
472 ssGBLUP than with PBLUP, in subsequent genetic evaluations when records on animals in the  
473 validation generation were excluded. However, when records on animals in the validation  
474 generation were included in subsequent genetic evaluations, validation accuracies were  
475 generally similar between corresponding ssGBLUP and PBLUP scenarios. The fact that  
476 heritabilities were all relatively high (ranging from 0.24 to 0.58, Tables 2 to 5) could, at least  
477 partly, explain the absence of significant differences between ssGBLUP and PBLUP  
478 evaluations when records on animals in the validation generation were included in the  
479 subsequent genetic evaluations. It is a common knowledge that the higher the heritability, the  
480 higher the importance of own record and the lesser the importance of genomic information in  
481 genetic evaluations (e.g. [13]).

482 In our previous study [3], we observed no absolute bias when ssGBLUP was used in subsequent  
483 genetic evaluations, irrespective preselection type or intensity. However, in [3], we found  
484 absolute bias to be increasing with intensity of preselection when we used PBLUP in subsequent  
485 genetic evaluations. Patry et al [1,6,7] also reported significant absolute bias when subsequent  
486 genetic evaluations of genomically preselected were done with PBLUP, except when some  
487 pseudo-phenotypic information on preculled animals was included in the subsequent PBLUP  
488 evaluations. As we did not include (pseudo) phenotypic information on preculled animals in

489 our subsequent PBLUP evaluations, we expected to find significant absolute bias, which would  
490 increase from reference to VGP to MGP scenarios. However, in the current study absolute bias  
491 remained largely absent across all the three scenarios of genomic preselection, irrespective of  
492 whether ssGBLUP or PBLUP was used.

493 In the absence of selection, the expectation of regression coefficient of PYD on (G)EBV - an  
494 indicator of dispersion bias - is 0.5, because PYD only represents half of the breeding value of  
495 the parent. However, when validation animals are not a representative sample of all animals in  
496 their age group, the expectation of the regression coefficient decreases, depending on how much  
497 the validation animals deviate from a random sample of animals in their age group [22,23]. In  
498 the data used in this study, average daily gain traits had heavier weights in the breeding goals  
499 of the two lines than backfat and loin depth, so we expected that our genomic preselection  
500 would have a smaller impact on the regression coefficients for backfat and loin depth than for  
501 the two average daily gain traits. We however did not observe smaller regression coefficients  
502 or regression coefficients further away from 0.5 for average daily gain traits than for backfat  
503 and loin depth, neither with ssGBLUP nor with PBLUP.

504 Regression coefficient of PYD on (G)EBV generally decreased with more genomic  
505 preselection, but were in most cases only marginally different from the expected value of 0.5.  
506 The decrease was more pronounced with PBLUP than with ssGBLUP. In many instances, the  
507 regression coefficients of reference scenarios with PBLUP were greater than 0.5, and they (the  
508 regression coefficients) became closer to 0.5 with more preselection. In our previous study with  
509 a simulated dataset [3], we found that regression coefficients of TBV on (G)EBV were bigger  
510 and closer to the expected value of 1 when ssGBLUP was used in the subsequent genetic  
511 evaluations compared to when PBLUP was used. In [3], we also found that the regression  
512 coefficient became smaller as preselection intensity increased when PBLUP was used, and

513 remained similar irrespective of preselection intensity when ssGBLUP was used. The generally  
514 similar regression coefficients across the genomic preselection scenarios with ssGBLUP in this  
515 study further confirms that ssGBLUP is indeed able to prevent most of the impact of  
516 preselection on subsequent genetic evaluations, as we previously reported in [3]. We have no  
517 explanation as to why regression coefficients from PBLUP were greater than the expected  
518 value, and also greater than their corresponding values from ssGBLUP. In conclusion, absolute  
519 bias remained largely absent across the three genomic preselection scenarios, while with more  
520 preselection validation accuracy only showed small decreases and hardly any dispersion bias  
521 was induced.

#### 522 **Comparison of results across the two lines**

523 Even in the dam-line where the original genomic preselection was more intense and ratio of  
524 males with records to females with records in any generation was about 20:80, we generally did  
525 not observe significantly greater biases with more genomic preselection. Although in both lines  
526 validation accuracy decreased with more genomic preselection for all traits and with or without  
527 records on animals in the validation generation, generally we did not find bigger decreases in  
528 the dam-line than in the sire-line. However, corresponding validation accuracies were always  
529 higher in the sire-line than in the dam-line, despite the corresponding estimated heritabilities  
530 being higher in the dam-line than in the sire-line for some traits. Corresponding regression  
531 coefficients of PYD on GEBV were also closer to the expected value of 0.5 in the sire-line than  
532 in the dam-line except for loin depth, where they were closer to 0.5 in the dam-line than in the  
533 sire-line. The observed higher accuracies and regression coefficients closer to the expected  
534 value in the sire-line than in the dam-line can most likely be explained by the higher  
535 phenotyping and genotyping rates in the sire-line than the dam-line (Table 1).

536

537 **Genotypes of preculled animals did not affect the subsequent ssGBLUP evaluations**

538 In the subsequent ssGBLUP evaluations without records on animals in the validation  
539 generation, results from corresponding reference and VGP scenarios were exactly the same, at  
540 least up to two decimal places (Tables 2 and 3). However, in terms of data content, reference  
541 scenarios contained genotypes of the animals preculled in the corresponding VGP scenarios, in  
542 addition to all the data contained in the corresponding VGP scenarios (Table 1). The fact that  
543 results from these two scenarios were the same means that genotypes of the preculled animals  
544 did not affect the reference scenarios. In this study, most (about 95%) of the validation animals  
545 and their parents had genotypes. This supports the conclusion from our previous study [4], that  
546 genotypes of preculled animals are only useful in subsequent ssGBLUP evaluations of their  
547 preselected sibs when their parents are not genotyped.

548 **Potential additional sources of bias in ssGBLUP from our data**

549 In practical datasets as used in this study, it is difficult to completely rule out some mistakes in  
550 pedigree recording and in genotyping. At our genomic data quality control stage, genotypes of  
551 a few thousand animals were discarded because the animals did not meet the genomic data  
552 quality standard (of being genotyped for at least 90% of the SNP). Genotyping mistakes could  
553 still not be completely ruled out in the genomic data that passed quality control. In Tables 2 to  
554 5, we saw that for some traits, heritabilities were different for different preselection scenarios,  
555 even though the animals in the base generation were the same. This implies that different  
556 subsets of the same data gave rise to different estimated (co)variance components in the base  
557 generation, and that it is likely that after some of the genomic preselection scenarios were  
558 implemented, the estimated (co)variance components were different from their true values, at  
559 least for some of the traits. While these are all potential additional sources of bias in ssGBLUP  
560 evaluations, they are difficult to avoid in practice [10]. However, in general, we can say that  
561 these potential additional sources of bias did not cause significant bias in our ssGBLUP

562 evaluations, as both absolute and dispersion biases were in most cases absent, and even when  
563 present they were only marginal.

## 564 **Conclusions**

565 When subsequent genetic evaluations of preselected animals are done with ssGBLUP, either  
566 with or without records on animals in the validation generation, realized accuracy reduces with  
567 genomic preselection in the validation generation, and even more with genomic preselection in  
568 multiple generations. On the other hand, absolute bias is largely absent, and dispersion bias  
569 only increases marginally with more genomic preselection in the current generation or in all  
570 generations. Impact of recent and/or historical genomic preselection is minimal on subsequent  
571 genetic evaluations of selection candidates, if these subsequent genetic evaluations are  
572 performed using ssGBLUP.

## 573 **Declarations**

### 574 **Ethical approval**

575 The data used for this study were collected as part of routine data recording in a commercial  
576 breeding program. Samples collected for DNA extraction were used for routine diagnostic  
577 purposes of the breeding program. Data recording and sample collection were conducted in line  
578 with local laws on protection of animals.

### 579 **Availability of data**

580 The data used in the present study were provided by Topigs Norsvin, and are not publicly  
581 accessible.

### 582 **Funding**

583 This study was financially supported by the Dutch Ministry of Economic Affairs (TKI Agri &

584 Food project 16022) and the Breed4Food partners Cobb Europe, CRV, Hendrix Genetics and  
585 Topigs Norsvin. The use of the HPC cluster was made possible by CAT-AgroFood (Shared  
586 Research Facilities Wageningen UR).

### 587 **Competing interests**

588 The authors declare that they have no competing interests.

### 589 **Authors' contributions**

590 All authors participated in the conception and the design of the study and of the analysis of the  
591 dataset. RB provided the dataset, IJ analysed the dataset and wrote the first draft of the  
592 manuscript, and the other authors revised the manuscript. All authors read and approved the  
593 final manuscript.

### 594 **Acknowledgements**

595 The authors thank Marco Bink and Katrijn Peeters from Hendrix Genetics, John Henshall from  
596 Cobb Europe, and Chris Schrooten and Gerben de Jong from CRV, for their inputs towards the  
597 design of this study.

## 598 **References**

- 599 1. Patry C, Ducrocq V. Evidence of biases in genetic evaluations due to genomic preselection  
600 in dairy cattle. *J Dairy Sci.* 2011;94:1011–20.
- 601 2. Masuda Y, VanRaden PM, Misztal I, Lawlor TJ. Differing genetic trend estimates from  
602 traditional and genomic evaluations of genotyped animals as evidence of preselection bias in  
603 us holsteins. *J Dairy Sci.* 2018;101:5194–206.
- 604 3. Jibrila I, Napel J, Vandenplas J, Veerkamp RF, Calus MPL. Investigating the impact of  
605 preselection on subsequent single-step genomic blup evaluation of preselected animals. *Genet*  
606 *Sel Evol.* 2020;52.
- 607 4. Jibrila I, Vandenplas J, ten Napel J, Veerkamp RF, Calus MPL. Avoiding preselection bias  
608 in subsequent single-step genomic blup evaluations of genomically preselected animals. *J Anim*  
609 *Breed Genet.* 2021;138: 432–41.

- 610 5. Shabalina T, Pimentel ECG, Edel C, Plieschke L, Emmerling R, Götz K-U. Short  
611 communication: the role of genotypes from animals without phenotypes in single-step genomic  
612 evaluations. *J Dairy Sci.* 2017;100:8277–81.
- 613 6. Patry C, Ducrocq V. Accounting for genomic pre-selection in national blup evaluations in  
614 dairy cattle. *Genet Sel Evol.* 2011;43.
- 615 7. Patry C, Jorjani H, Ducrocq V. Effects of a national genomic preselection on the international  
616 genetic evaluations. *J Dairy Sci.* 2013;96:3272–84.
- 617 8. Henderson CR. Best linear unbiased estimation and prediction under a selection model.  
618 *Biometris.* 1975;31:423–47.
- 619 9. Pollak EJ, van der Werf J, Quaas RL. Selection bias and multiple trait evaluation. *J Dairy*  
620 *Sci.* 1984;67:1590–5.
- 621 10. Tsuruta S, Lourenco DAL, Masuda Y, Misztal I, Lawlor TJ. Controlling bias in genomic  
622 breeding values for young genotyped bulls. *J Dairy Sci.* 2019;102:9956–70.
- 623 11. Vitezica ZG, Aguilar I, Misztal I, Legarra A. Bias in genomic predictions for populations  
624 under selection. *Genet Res (Camb).* 2011;93:357–66.
- 625 12. Hsu W-L, Garrick DJ, Fernando RL. The accuracy and bias of single-step genomic  
626 prediction for populations under selection. *Genes|Genomes|Genetics.* 2017;7:2685–94.
- 627 13. Mrode RA. Linear models for the prediction of animal breeding values. 3rd ed. 2014.
- 628 14. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MAR, Bender D, et al. PLINK: a  
629 tool set for whole-genome association and population-based linkage analyses. *Am J Hum*  
630 *Genet.* 2007;81:559–75.
- 631 15. Gilmour AR, Gogel BJ, Cullis BR, Thompson R. ASReml user guide release 3.0. VSN Int.  
632 Ltd. 2009. p. 275.
- 633 16. ten Napel J, Vandenplas J, Lidauer M, Strandén I, Taskinen M, Mäntysaari E, et al.  
634 MiXB LUP: a user-friendly software for large genetic evaluation systems. 2020. p. 62.
- 635 17. Aguilar I, Misztal I, Johnson DL, Legarra A, Tsuruta S, Lawlor TJ. Hot topic: a unified  
636 approach to utilize phenotypic, full pedigree, and genomic information for genetic evaluation  
637 of holstein final score. *J Dairy Sci.* 2010;93:743–52.
- 638 18. Christensen OF, Lund MS. Genomic prediction when some animals are not genotyped.  
639 *Genet Sel Evol.* 2010;42.
- 640 19. VanRaden PM. Efficient methods to compute genomic predictions. *J Dairy Sci.*  
641 2008;91:4414–23.
- 642 20. Cameron ND. Selection indices and prediction of genetic merit in animal breeding. 1997.
- 643 21. Duenk P, Calus MPL, Wientjes YCJ, Breen VP, Henshall JM, Hawken R, et al. Validation

644 of genomic predictions for body weight in broilers using crossbred information and considering  
645 breed-of-origin of alleles. Genet Sel Evol. 2019;51.

646 22. Mäntysaari EA, Liu Z, VanRaden P. Interbull validation test for genomic evaluations.  
647 Interbull Bull. 2010;17.

648 23. Mäntysaari EA, Koivula M. GEBV validation test revisited. Interbull Bull. 2012;45.

## 649 **Figures**

### 650 **Figure 1 Schematic representation of the animals included in the subsequent genetic** 651 **evaluations following each genomic preselection scenario**

652 Following the reference scenario, all animals in the figure were included in the subsequent  
653 evaluations. In the VGP scenario, only the culled animals in the validation generation were  
654 excluded from the subsequent evaluations. Finally, in the MGP scenario, all culled animals in  
655 all generations were excluded from the subsequent evaluations. Selection and culling here refer  
656 to those conducted by Topigs Norsvin as part of the company's routine practices.

## 657 **Additional files**

### 658 **Additional file 1 Table S1**

659 Format: .docx

660 Title: Estimated additive genetic and residual variances in the sire-line

661 Description: The additive genetic and residual variances that resulted to different heritability  
662 estimates for the same traits under different scenarios of subsequent genetic evaluations, in the  
663 sire line

### 664 **Additional file 2 Table S2**

665 Format: .docx

666 Title: Estimated additive genetic and residual variances in the dam-line

667 Description: The additive genetic and residual variances that resulted to different heritability  
668 estimates for the same traits under different scenarios of subsequent genetic evaluations, in the  
669 dam line

# Figures

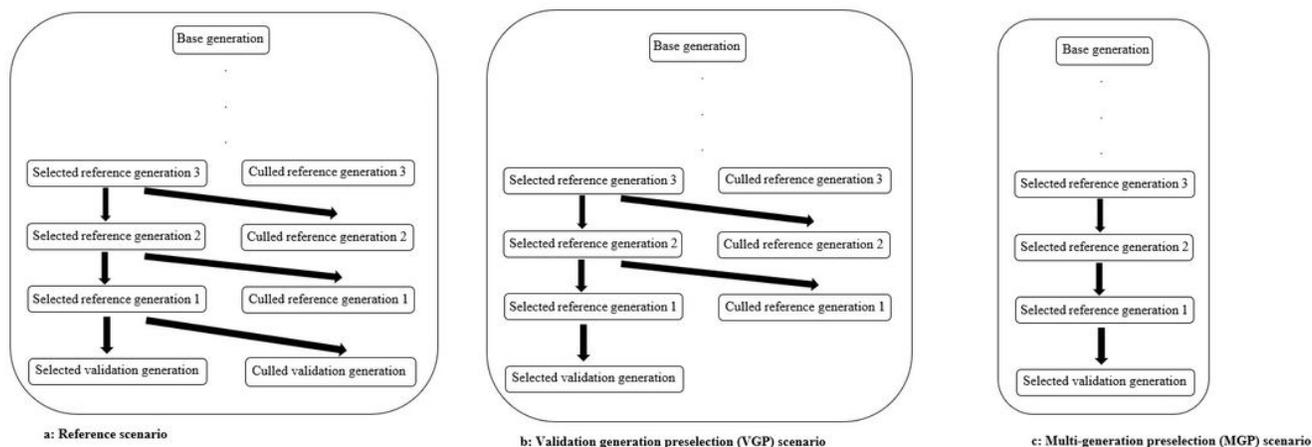


Figure 1

Schematic representation of the animals included in the subsequent genetic evaluations following each genomic preselection scenario. Following the reference scenario, all animals in the figure were included in the subsequent evaluations. In the VGP scenario, only the culled animals in the validation generation were excluded from the subsequent evaluations. Finally, in the MGP scenario, all culled animals in all generations were excluded from the subsequent evaluations. Selection and culling here refer to those conducted by Topigs Norsvin as part of the company's routine practices.

## Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [Additionalfile1.docx](#)
- [Additionalfile2.docx](#)