

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>*  
4 *and 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](https://orcid.org/0000-0002-5683-1263)

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

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

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

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

16 MPLC: [mario.calus@wur.nl](mailto:mario.calus@wur.nl); ORCID: [0000-0002-3213-704X](https://orcid.org/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.  
29 However, preselection almost always takes place in animal breeding programs, so it is  
30 difficult, if not impossible, to have a dataset without preselection. Hence most studies on  
31 preselection used simulated datasets, concluding that subsequent genomic estimated breeding  
32 values (GEBV) from single-step genomic best linear unbiased prediction (ssGBLUP) are  
33 unbiased. The aim of this study was to investigate the impact of genomic preselection, using  
34 real data, on accuracy 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.  
39 Per 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  
42 validation generation, or ii) in all generations. In computing accuracy and bias, we compared  
43 GEBV 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  
47 animals had records or not. Bias too was largely absent, and did not differ greatly among  
48 corresponding 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  
60 the initial stages of this selection are called preselection [1–3]. Selection candidates that  
61 survive preselection are called preselected animals [1–3], and those that do not are called  
62 preculled animals [3,4]. Preselection aims to reduce costs and efforts spent on animals that are  
63 not 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  
68 preselected animals. Properly assessing the impact of preselection on subsequent genetic  
69 evaluation of preselected animals requires a scenario without preselection, against which  
70 scenarios with preselection can be compared. Because in animal breeding programmes  
71 preselection almost always takes place, it is difficult, if not impossible, to have a scenario  
72 without preselection. This is why most studies available on preselection used simulated  
73 datasets [e.g. 1–5]. Those studies have shown that when a subsequent genetic evaluation of  
74 preselected animals is done using pedigree-based best linear unbiased prediction (PBLUP),  
75 genomic preselection results in accuracy loss and bias in the estimated breeding values (EBV)  
76 of preselected animals [1,6–9]. Some of these studies [6–9] further showed that the accuracy  
77 loss and bias caused by genomic preselection can be avoided if the information on preculled  
78 animals that was utilized at preselection is included in the subsequent PBLUP evaluation. On  
79 the other hand, our previous works [3,4] have shown that when the subsequent genetic  
80 evaluation is done with single-step genomic BLUP (ssGBLUP), genomic EBV (GEBV) of  
81 preselected animals are estimated without bias. We [4] further showed that to avoid genomic

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

## 100 **Methods**

### 101 **Data**

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

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

118 We obtained pig production traits data on one sire-line and one dam-line from Topigs  
119 Norsvin. These data were collected between 1970 and 2020, and the traits were average daily  
120 gain during performance testing, average daily gain throughout the lifetime, backfat, and loin  
121 depth. Topigs Norsvin (pre)selected both lines on these production traits. However, there was  
122 more emphasis on reproduction traits than on production traits in the dam-line. Details on the  
123 amount of data utilized in this study are in Table 1. The data were recorded on originally  
124 preselected animals (i.e. the animals preselected by Topigs Norsvin), with the sire-line being  
125 much more balanced than the dam-line, in terms of proportions males and females with  
126 records per generation (ratio of males with records to females with records is about 50:50 in  
127 the sire-line and about 20:80 in the dam-line). We studied impact of genomic preselection in  
128 the two lines separately, because the traits we studied had different weights in breeding goals  
129 of the two lines. The cut-off date to split the data into reference and validation populations  
130 was 31<sup>st</sup> January, 2017 for the sire line, and 31<sup>st</sup> December, 2015 for the dam-line. In the

131 pedigree, animals with one or both parents missing were assigned to genetic groups,  
132 according to line and year of birth of each animal.

### 133 **Genomic data and quality control**

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

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

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

$$152 \quad \mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{W}\mathbf{p} + \mathbf{Z}\mathbf{u} + \mathbf{e}, \quad (\text{eq. 1}),$$

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

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

$$157 \quad y_{ci} = \hat{u}_i + \hat{e}_i \quad (\text{eq. 2}).$$

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

### 161 **Preselection**

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

### 176 **Subsequent genetic evaluations**

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

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

$$189 \quad \mathbf{y} = \mathbf{x}\mathbf{b} + \mathbf{Z}\mathbf{u} + \mathbf{e} \quad (\text{eq. 3}),$$

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

195 **Figure 1 here**

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



208

209

210 **Table 1** Data utilized in subsequent ssGBLUP<sup>a</sup> evaluations following each preselection  
 211 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

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

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

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

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

## 220 Implementation of single-step GBLUP

221 The inverse of the combined pedigree-genomic relationship ( $\mathbf{H}^{-1}$ ) was obtained as follows  
 222 [17,18]:

$$223 \quad \mathbf{H}^{-1} = \mathbf{A}^{-1} + \begin{bmatrix} 0 & 0 \\ 0 & (0.95\mathbf{G}_t + 0.05\mathbf{A}_{22})^{-1} - \mathbf{A}_{22}^{-1} \end{bmatrix} \quad (\text{eq. 4}),$$

224 where  $\mathbf{A}^{-1}$  was the inverse of the pedigree relationship matrix, and  $\mathbf{A}_{22}$  was part of the  
 225 pedigree relationship matrix referring to genotyped animals. We considered inbreeding in

226 setting up both  $\mathbf{A}^{-1}$  and  $\mathbf{A}_{22}$  to avoid bias caused by ignoring inbreeding (Tsuruta et al.,  
227 2019). The genomic relationship matrix  $\mathbf{G}_t$  was computed as follows:

$$228 \quad \mathbf{G}_t = (1 - \bar{f}_p)\mathbf{G}_r + 2\bar{f}_p \mathbf{1}\mathbf{1}' \quad (\text{eq. 5}),$$

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

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

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

$$247 \quad PYD_i = \frac{\sum_{p=1}^n y_{cp} - g_m}{n} \quad (\text{eq. 6}),$$

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

257 
$$\frac{1/4nh^2}{1+1/4(n-1)h^2} \quad (\text{eq. 7}),$$

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

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

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

## 281 **Results**

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

288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301

302  
303  
304  
305  
306  
307  
308  
309  
310

**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

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

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

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

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

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

321  
322  
323  
324  
325  
326  
327  
328

329  
330  
331  
332  
333  
334  
335  
336  
337

**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

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

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

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

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

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

348  
349

350  
351  
352  
353  
354  
355  
356  
357  
358

**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

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

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

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

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

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

369  
370

371  
372  
373  
374  
375  
376  
377  
378  
379

**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

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

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

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

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

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

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



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

412 **Subsequent ssGBLUP evaluations without records on animals in the validation**  
413 **generation**

414 Without records on animals in the validation generation in the subsequent ssGBLUP  
415 evaluations, all results for the reference and VGP scenarios were the same. Just like when

416 records on animals in the validation generation were included, here too, estimated heritability  
417 increased from reference and VGP to MGP scenarios, and in this case for all traits in both  
418 lines. Validation accuracy also decreased from reference and VGP to MGP scenarios, and in  
419 this case with bigger decreases compared to when records on animals in the validation  
420 generation were included. Absolute bias was also largely absent for backfat and loin depth for  
421 both lines, and showed no particular pattern for average daily gain traits for the two lines.  
422 Even for the average daily gain traits, it was still small, with  $\pm 0.14$  additive genetic SD being  
423 the highest value (Tables 2 and 3). Regression coefficients of PYD on GEBV were similar to  
424 their corresponding value when records on animals in the validation generation were included.  
425 The only exception were all scenarios for backfat in the sire-line, where the regression  
426 coefficients of PYD on GEBV appeared to be lower than their corresponding values when  
427 records on animals in the validation generation were included. For both lines, the regression  
428 coefficients ranged from 0.41 (for the MGP scenario for backfat in the dam-line) to 0.53 (for  
429 the reference and VGP scenarios for loin depth in the dam-line).

#### 430 **Subsequent genetic evaluations with PBLUP**

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

## 440 **Discussion**

441 In this study, we investigated the impact of genomic preselection on subsequent ssGBLUP  
442 evaluations of preselected animals, using real data from an ongoing pig breeding program in  
443 which preselection has taken place, by retrospectively implementing additional layers of  
444 preselection. The data was on production traits of pigs from one sire-line and one dam-line.  
445 Per line, we implemented three genomic preselection scenarios. We used pre-corrected  
446 phenotypes as records in the subsequent genetic evaluations, and progeny yield deviation  
447 (PYD) as the proxy for TBV. We did the subsequent genetic evaluations either with or  
448 without records on animals in the validation generation, and in all cases without progeny of  
449 validation animals. In both lines, for all traits and with or without records on validation  
450 animals, absolute bias was largely absent across the three genomic preselection scenarios,  
451 while with more preselection validation accuracy only showed small decreases and hardly any  
452 dispersion bias was induced.

453 In the two scenarios with additional genomic preselection (i.e. VGP and MGP scenarios), the  
454 preselected animals in every generation were the animals that in reality were selected and  
455 produced progeny in the next generation, and the preculled animals were those animals that  
456 were in reality culled after performance testing. Thus, these two scenarios represent either i)  
457 situations in which all the selection in a generation is done in only one stage, after selection  
458 candidates have own records, or ii) situations in which an additional selection stage is  
459 implemented after preselected animals have had progeny. While neither of these cases is true  
460 for the data we used, the scenarios we implemented enabled us to investigate the impact of  
461 genomic preselection on subsequent genetic evaluations of preselected animals using real  
462 data, by including different amounts of pedigree, genomic and phenotypic information in the  
463 subsequent genetic evaluations we implemented. The validation accuracy we computed as the

464 correlation between (G)EBV and PYD is not numerically the same as the accuracy of  
465 predicting TBV, since variance of PYD has some non-genetic component, in addition to  
466 genetic component [13]. However, the two accuracies are proportional to each other, and this  
467 enabled us to make comparison among subsequent genetic evaluation scenarios [21].

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

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

489 In our previous study [3], we observed no absolute bias when ssGBLUP was used in  
490 subsequent genetic evaluations, irrespective of preselection type or intensity. However, in [3],  
491 we found absolute bias to be increasing with intensity of preselection when we used PBLUP  
492 in subsequent genetic evaluations. Patry et al [1,6,7] also reported significant absolute bias  
493 when subsequent genetic evaluations of genomically preselected were done with PBLUP,  
494 except when some pseudo-phenotypic information on preculled animals was included in the  
495 subsequent PBLUP evaluations. As we did not include (pseudo) phenotypic information on  
496 preculled animals in our subsequent PBLUP evaluations, we expected to find significant  
497 absolute bias, which would increase from reference to VGP to MGP scenarios. However, in  
498 the current study absolute bias remained largely absent across all the three scenarios of  
499 genomic preselection, irrespective of whether ssGBLUP or PBLUP was used.

500 In the absence of selection, the expectation of regression coefficient of PYD on (G)EBV - an  
501 indicator of dispersion bias - is 0.5, because PYD only represents half of the breeding value of  
502 the parent. However, when validation animals are not a representative sample of all animals in  
503 their age group, the expectation of the regression coefficient decreases, depending on how  
504 much the validation animals deviate from a random sample of animals in their age group  
505 [22,23]. In the data used in this study, average daily gain traits had heavier weights in the  
506 breeding goals of the two lines than backfat and loin depth, so we expected that our genomic  
507 preselection would have a smaller impact on the regression coefficients for backfat and loin  
508 depth than for the two average daily gain traits. We however did not observe smaller  
509 regression coefficients or regression coefficients further away from 0.5 for average daily gain  
510 traits than for backfat and loin depth, neither with ssGBLUP nor with PBLUP.

511 Regression coefficient of PYD on (G)EBV generally decreased with more genomic  
512 preselection, but were in most cases only marginally different from the expected value of 0.5.

513 The decrease was more pronounced with PBLUP than with ssGBLUP. In many instances, the  
514 regression coefficients of reference scenarios with PBLUP were greater than 0.5, and they  
515 (the regression coefficients) became closer to 0.5 with more preselection. In our previous  
516 study with a simulated dataset [3], we found that regression coefficients of TBV on (G)EBV  
517 were bigger and closer to the expected value of 1 when ssGBLUP was used in the subsequent  
518 genetic evaluations compared to when PBLUP was used. In [3], we also found that the  
519 regression coefficient became smaller as preselection intensity increased when PBLUP was  
520 used, and remained similar irrespective of preselection intensity when ssGBLUP was used.  
521 The generally similar regression coefficients across the genomic preselection scenarios with  
522 ssGBLUP in this study further confirms that ssGBLUP is indeed able to prevent most of the  
523 impact of preselection on subsequent genetic evaluations, as we previously reported in [3].  
524 We have no explanation as to why regression coefficients from PBLUP were greater than the  
525 expected value, and also greater than their corresponding values from ssGBLUP. In  
526 conclusion, absolute bias remained largely absent across the three genomic preselection  
527 scenarios, while with more preselection validation accuracy only showed small decreases and  
528 hardly any dispersion bias was induced.

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

530 Even in the dam-line where the original genomic preselection was more intense and ratio of  
531 males with records to females with records in any generation was about 20:80, we generally  
532 did not observe significantly greater biases with more genomic preselection. Although in both  
533 lines validation accuracy decreased with more genomic preselection for all traits and with or  
534 without records on animals in the validation generation, generally we did not find bigger  
535 decreases in the dam-line than in the sire-line. However, corresponding validation accuracies  
536 were always higher in the sire-line than in the dam-line, despite the corresponding estimated  
537 heritabilities being higher in the dam-line than in the sire-line for some traits. Corresponding

538 regression coefficients of PYD on GEBV were also closer to the expected value of 0.5 in the  
539 sire-line than in the dam-line except for loin depth, where they were closer to 0.5 in the dam-  
540 line than in the sire-line. The observed higher accuracies and regression coefficients closer to  
541 the expected value in the sire-line than in the dam-line can most likely be explained by the  
542 higher phenotyping and genotyping rates in the sire-line than the dam-line (Table 1).

543

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

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

#### 555 **Potential additional sources of bias in ssGBLUP from our data**

556 In practical datasets as used in this study, it is difficult to completely rule out some mistakes  
557 in pedigree recording and in genotyping. At our genomic data quality control stage, genotypes  
558 of a few thousand animals were discarded because the animals did not meet the genomic data  
559 quality standard (of being genotyped for at least 90% of the SNP). Genotyping mistakes could  
560 still not be completely ruled out in the genomic data that passed quality control. In Tables 2 to  
561 5, we saw that for some traits, heritabilities were different for different preselection scenarios,

562 even though the animals in the base generation were the same. This implies that different  
563 subsets of the same data gave rise to different estimated (co)variance components in the base  
564 generation, and that it is likely that after some of the genomic preselection scenarios were  
565 implemented, the estimated (co)variance components were different from their true values, at  
566 least for some of the traits. While these are all potential additional sources of bias in  
567 ssGBLUP evaluations, they are difficult to avoid in practice [10]. However, in general, we  
568 can say that these potential additional sources of bias did not cause significant bias in our  
569 ssGBLUP evaluations, as both absolute and dispersion biases were in most cases absent, and  
570 even when present they were only marginal.

## 571 **Conclusions**

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

## 580 **Declarations**

### 581 **Ethical approval**

582 The data used for this study were collected as part of routine data recording in a commercial  
583 breeding program. Samples collected for DNA extraction were used for routine diagnostic  
584 purposes of the breeding program. Data recording and sample collection were conducted in



585 line with local laws on protection of animals.

#### 586 **Availability of data**

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

#### 589 **Funding**

590 This study was financially supported by the Dutch Ministry of Economic Affairs (TKI Agri &  
591 Food project 16022) and the Breed4Food partners Cobb Europe, CRV, Hendrix Genetics and  
592 Topigs Norsvin. The use of the HPC cluster was made possible by CAT-AgroFood (Shared  
593 Research Facilities Wageningen UR).

#### 594 **Competing interests**

595 The authors declare that they have no competing interests.

#### 596 **Authors' contributions**

597 All authors participated in the conception and the design of the study and of the analysis of  
598 the dataset. RB provided the dataset, IJ analysed the dataset and wrote the first draft of the  
599 manuscript, and the other authors revised the manuscript. All authors read and approved the  
600 final manuscript.

#### 601 **Acknowledgements**

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

#### 605 **References**

606 1. Patry C, Ducrocq V. Evidence of biases in genetic evaluations due to genomic preselection

- 607 in dairy cattle. *J Dairy Sci.* 2011;94:1011–20.
- 608 2. Masuda Y, VanRaden PM, Misztal I, Lawlor TJ. Differing genetic trend estimates from  
609 traditional and genomic evaluations of genotyped animals as evidence of preselection bias in  
610 us holsteins. *J Dairy Sci.* 2018;101:5194–206.
- 611 3. Jibrila I, Napel J, Vandenplas J, Veerkamp RF, Calus MPL. Investigating the impact of  
612 preselection on subsequent single-step genomic blup evaluation of preselected animals.  
613 *Genet Sel Evol.* 2020;52.
- 614 4. Jibrila I, Vandenplas J, ten Napel J, Veerkamp RF, Calus MPL. Avoiding preselection bias  
615 in subsequent single-step genomic blup evaluations of genomically preselected animals. *J*  
616 *Anim Breed Genet.* 2021;138: 432–41.
- 617 5. Shabalina T, Pimentel ECG, Edel C, Plieschke L, Emmerling R, Götz K-U. Short  
618 communication: the role of genotypes from animals without phenotypes in single-step  
619 genomic evaluations. *J Dairy Sci.* 2017;100:8277–81.
- 620 6. Patry C, Ducrocq V. Accounting for genomic pre-selection in national blup evaluations in  
621 dairy cattle. *Genet Sel Evol.* 2011;43.
- 622 7. Patry C, Jorjani H, Ducrocq V. Effects of a national genomic preselection on the  
623 international genetic evaluations. *J Dairy Sci.* 2013;96:3272–84.
- 624 8. Henderson CR. Best linear unbiased estimation and prediction under a selection model.  
625 *Biometris.* 1975;31:423–47.
- 626 9. Pollak EJ, van der Werf J, Quaas RL. Selection bias and multiple trait evaluation. *J Dairy*  
627 *Sci.* 1984;67:1590–5.
- 628 10. Tsuruta S, Lourenco DAL, Masuda Y, Misztal I, Lawlor TJ. Controlling bias in genomic  
629 breeding values for young genotyped bulls. *J Dairy Sci.* 2019;102:9956–70.
- 630 11. Vitezica ZG, Aguilar I, Misztal I, Legarra A. Bias in genomic predictions for populations  
631 under selection. *Genet Res (Camb).* 2011;93:357–66.
- 632 12. Hsu W-L, Garrick DJ, Fernando RL. The accuracy and bias of single-step genomic  
633 prediction for populations under selection. *Genes|Genomes|Genetics.* 2017;7:2685–94.
- 634 13. Mrode RA. Linear models for the prediction of animal breeding values. 3rd ed. 2014.
- 635 14. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MAR, Bender D, et al. PLINK: a  
636 tool set for whole-genome association and population-based linkage analyses. *Am J Hum*  
637 *Genet.* 2007;81:559–75.
- 638 15. Gilmour AR, Gogel BJ, Cullis BR, Thompson R. ASReml user guide release 3.0. VSN  
639 Int. Ltd. 2009. p. 275.
- 640 16. ten Napel J, Vandenplas J, Lidauer M, Strandén I, Taskinen M, Mäntysaari E, et al.  
641 MiXBLUP: a user-friendly software for large genetic evaluation systems. 2020. p. 62.

- 642 17. Aguilar I, Misztal I, Johnson DL, Legarra A, Tsuruta S, Lawlor TJ. Hot topic: a unified  
643 approach to utilize phenotypic, full pedigree, and genomic information for genetic evaluation  
644 of holstein final score. *J Dairy Sci.* 2010;93:743–52.
- 645 18. Christensen OF, Lund MS. Genomic prediction when some animals are not genotyped.  
646 *Genet Sel Evol.* 2010;42.
- 647 19. VanRaden PM. Efficient methods to compute genomic predictions. *J Dairy Sci.*  
648 2008;91:4414–23.
- 649 20. Cameron ND. Selection indices and prediction of genetic merit in animal breeding. 1997.
- 650 21. Duenk P, Calus MPL, Wientjes YCJ, Breen VP, Henshall JM, Hawken R, et al.  
651 Validation of genomic predictions for body weight in broilers using crossbred information  
652 and considering breed-of-origin of alleles. *Genet Sel Evol.* 2019;51.
- 653 22. Mäntysaari EA, Liu Z, VanRaden P. Interbull validation test for genomic evaluations.  
654 *Interbull Bull.* 2010;17.
- 655 23. Mäntysaari EA, Koivula M. GEBV validation test revisited. *Interbull Bull.* 2012;45.

## 656 **Figures**

### 657 **Figure 1 Schematic representation of the animals included in the subsequent genetic** 658 **evaluations following each genomic preselection scenario**

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

## 664 **Additional files**

### 665 **Additional file 1 Table S1**

666 Format: .docx

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

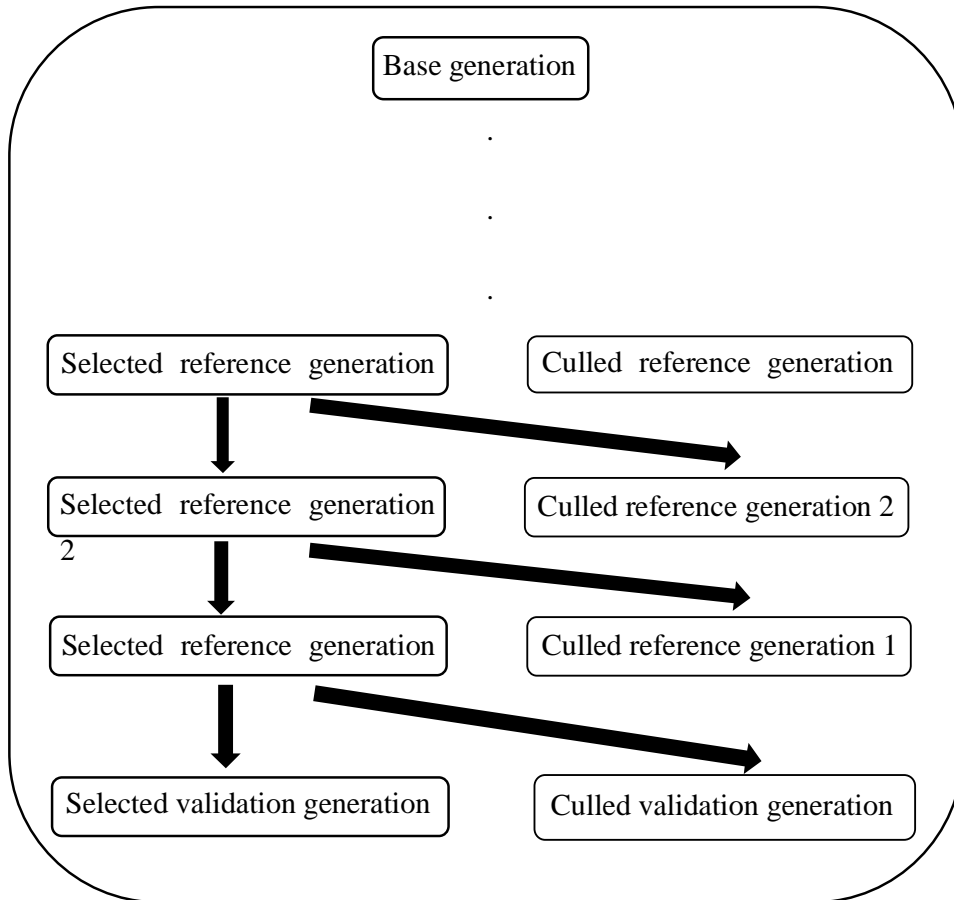
668 Description: The additive genetic and residual variances that resulted to different heritability  
669 estimates for the same traits under different scenarios of subsequent genetic evaluations, in  
670 the sire line

### 671 **Additional file 2 Table S2**

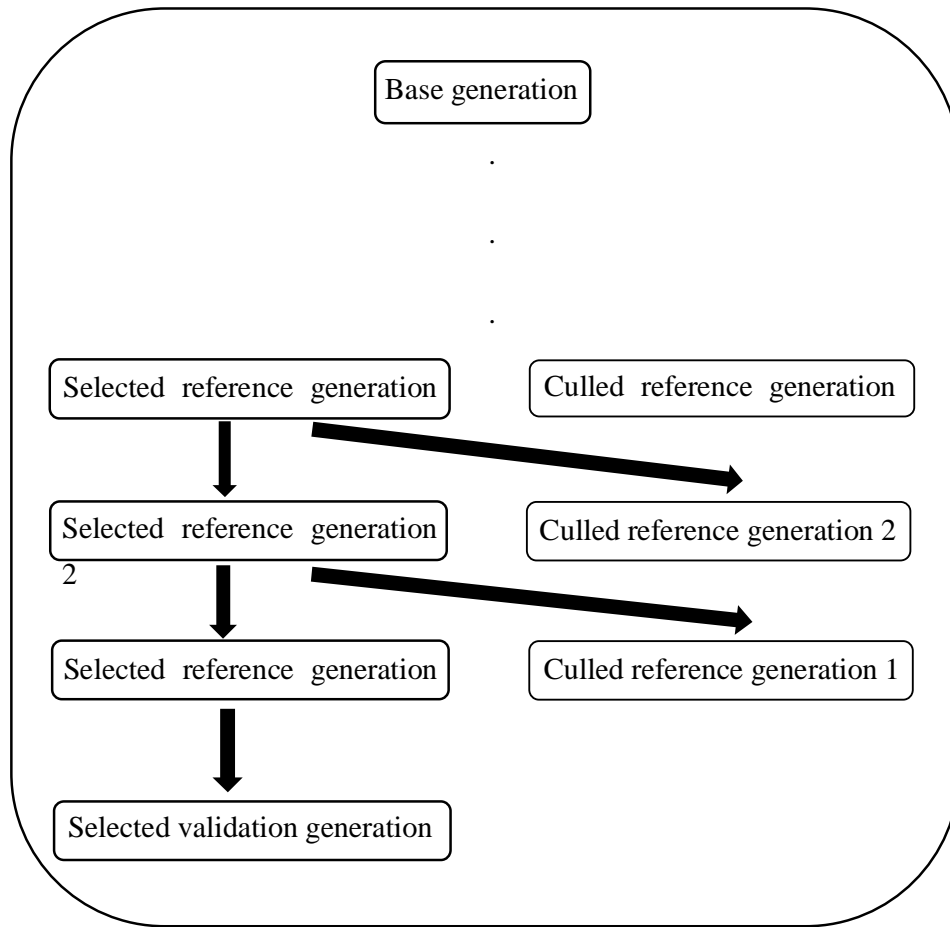
672 Format: .docx

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

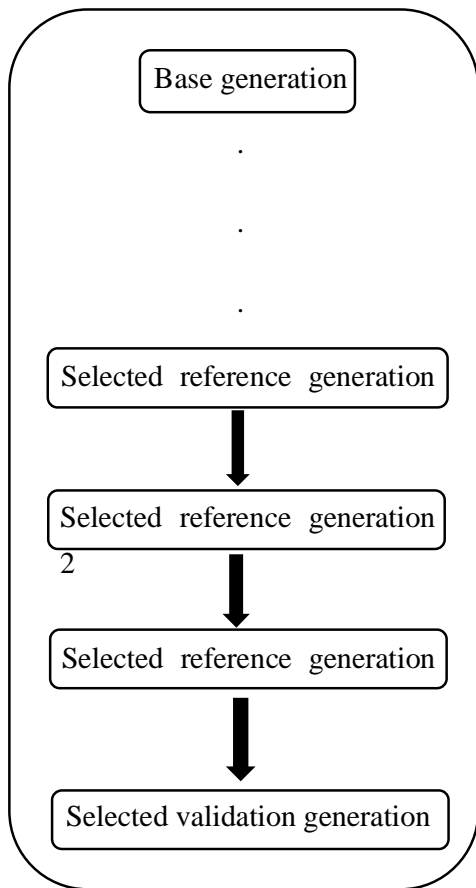
674 Description: The additive genetic and residual variances that resulted to different heritability  
675 estimates for the same traits under different scenarios of subsequent genetic evaluations, in  
676 the dam line



**a: Reference scenario**



**b: Validation generation preselection (VGP) scenario**



**c: Multi-generation preselection (MGP) scenario**