

1 **Running head:** REML with the algorithm for proven and young

2

3 **Is single-step genomic REML with the algorithm for proven and young more**
4 **computationally efficient when less generations of data are present?**

5

6 Vinícius Silva Junqueira*^{1,2}, Daniela Lourenco†, Yutaka Masuda†, Fernando Flores Cardoso‡,
7 Paulo Sávio Lopes*, Fabyano Fonseca e Silva*, Ignacy Misztal†

8

9 ² Breeding Research Department, Bayer Crop Science, Uberlândia, Minas Gerais, Brazil

10 * Departamento de Zootecnia, Universidade Federal de Viçosa, Viçosa, Minas Gerais, Brazil

11 † Department of Dairy and Animal Science, University of Georgia, Athens, Georgia, United
12 States

13 ‡ Empresa Brasileira de Pesquisa Agropecuária (EMBRAPA) Pecuária Sul, Bagé, Rio Grande do
14 Sul, Brasil

15

16

17

18 ¹ Corresponding author: viniciussilva.junqueira@bayer.com

19

20 **Lay Summary**

21 The estimation of variance components is computationally expensive under large-scale
22 genetic evaluations due to several inversions of the coefficient matrix. Variance components are
23 used as parameters for estimating breeding values in mixed model equations (MME). However,
24 resulting breeding values are not Best Linear Unbiased Predictions (BLUP) unless the variance
25 components approach the true parameters. The increasing availability of genomic data requires
26 the development of new methods for improving the efficiency of variance component
27 estimations. Therefore, this study aimed to reduce the costs of single-step genomic REML
28 (ssGREML) with the Algorithm for Proven and Young (APY) for estimating variance
29 components with truncated pedigree and phenotypes. In addition, we investigated the influence
30 of truncation on variance components and genetic parameter estimates. Under APY, the size of
31 the core group influences the similarity of breeding values and their reliability compared to the
32 full genomic matrix. In this study, we found that to ensure reliable variance component
33 estimation it is required to consider a core size that corresponds to the number of largest
34 eigenvalues explaining around 98% of the total variation in \mathbf{G} to avoid biased parameters. In
35 terms of costs, the use of APY slightly decreased the time for ordering and symbolic
36 factorization with no impact on estimations.

37

38 **Teaser Text**

39 Estimation of variance components is becoming computationally challenging due to the
40 increasing size of genomic information. We investigated the impacts of using the algorithm for
41 proven and young (APY) in genetic evaluations. The use of APY has no impact on variance
42 components and genetic parameters estimation.

43 **Abstract:**

44 Efficient computing techniques allow the estimation of variance components for virtually any
45 traditional dataset. When genomic information is available, variance components can be
46 estimated using genomic REML (GREML). If only a portion of the animals have genotypes,
47 single-step GREML (ssGREML) is the method of choice. The genomic relationship matrix (\mathbf{G})
48 used in both cases is dense, limiting computations depending on the number of genotyped
49 animals. The algorithm for proven and young (APY) can be used to create a sparse inverse of \mathbf{G}
50 (\mathbf{G}_{APY}^{-1}) with close to linear memory and computing requirements. In ssGREML, the inverse of
51 the realized relationship matrix (\mathbf{H}^{-1}) also includes the inverse of the pedigree relationship
52 matrix, which can be dense with long pedigree, but sparser with short. The main purpose of this
53 study was to investigate whether costs of ssGREML can be reduced using APY with truncated
54 pedigree and phenotypes. We also investigated the impact of truncation on variance components
55 estimation when different numbers of core animals are used in APY. Simulations included 150K
56 animals from 10 generations, with selection. Phenotypes ($h^2 = 0.3$) were available for all animals
57 in generations 1-9. A total of 30K animals in generations 8 and 9, and 15K validation animals in
58 generation 10 were genotyped for 52,890 SNP. Average information REML and ssGREML with
59 \mathbf{G}^{-1} and \mathbf{G}_{APY}^{-1} using 1K, 5K, 9K, and 14K core animals were compared. Variance components
60 are impacted when the core group in APY represents the number of eigenvalues explaining a

61 small fraction of the total variation in **G**. The most time-consuming operation was the inversion,
62 with more than 50% of the total time. Next, numerical factorization consumed nearly 30% of the
63 total computing time. On average, a 7% decrease in the computing time for ordering was
64 observed by removing each generation of data. APY can be successfully applied to create the
65 inverse of the genomic relationship matrix used in ssGREML for estimating variance
66 components. To ensure reliable variance component estimation, it is important to use a core size
67 that corresponds to the number of largest eigenvalues explaining around 98% of total variation in
68 **G**. When APY is used, pedigrees can be truncated to increase the sparsity of **H** and slightly
69 reduce computing time for ordering and symbolic factorization, with no impact on the estimates.

70 **Keywords:** variance components, genomic information, sparse genomic matrix, old data

71 **Abbreviations**

72	A	pedigree relationship matrix
73	AIREML	average information restricted maximum likelihood
74	APY	algorithm for proven and young
75	BLUP	best linear unbiased prediction
76	EBV	estimated breeding value
77	G	genomic matrix
78	G_{APY}	genomic matrix created using APY
79	GEHV	genomic enhanced breeding value
80	GREML	genomic restricted maximum likelihood
81	IOD	iteration on data
82	LHS	left hand side of mixed model equations

83	MME	mixed model equations
84	QTL	quantitative trait loci
85	REML	restricted maximum likelihood
86	ssGBLUP	single step genomic BLUP
87	ssGREML	single step genomic restricted maximum likelihood
88	YAMS	yet another MME solver

89

90

Introduction

91 Restricted maximum likelihood (REML), described by Patterson and Thompson (1971),
92 is a popular method for parameter estimation. Because it uses the mixed model equations
93 (Henderson, 1975), it is resistant to selection bias, and efficient implementations are currently
94 available. With the Average Information (AI) algorithm, convergence is often achieved in a few
95 rounds. With traces obtained by sparse matrix factorization and inversion (Meyer, 1997),
96 computing variance components is feasible even with large models.

97 When genomic information is available, two versions of REML may be applicable. When
98 only genotyped animals have phenotypes, genomic REML (GREML) can be applied with a
99 genomic relationship matrix (**G**). In general, such a matrix is dense, and the cost of dense matrix
100 operations would limit computations depending on the models. When only a fraction of animals
101 are genotyped, a single-step genomic REML is applicable (ssGREML). In the latter, the
102 combined relationship matrix (**H**) has dense blocks due to the genomic information, limiting the
103 efficiency of sparse matrix operations. Lately, Masuda et al. (2015) developed a sparse matrix
104 package YAMS that identifies dense blocks and computes them efficiently. For ssGREML, with

105 genomic computation, such a package resulted in up to 100 times speedup, allowing four trait
106 models with 20,000 genotyped animals (Masuda et al., 2015).

107 In general, it is of interest to include many genotyped animals in parameter estimation
108 and evaluations to account for genomic selection or pre-selection (Patry and Ducrocq, 2011). For
109 instance, the greatest reliability in a single-step genomic BLUP was obtained using 50% of the
110 heritability computed with a non-genomic REML (Misztal et al., 2017). The number of
111 genotyped animals is increasing fast for some species. As an example, almost 3 million Holsteins
112 have been genotyped in the US (https://queries.uscdcb.com/Genotype/cur_freq.html). However,
113 the cost of dense matrix operations with \mathbf{G} in REML using YAMS is quadratic for memory and
114 cubic for operations, which limits computations to around 50,000 animals.

115 The genomic information has a limited dimensionality due to the limited effective
116 population size (Stam, 1980; VanRaden, 2008; Misztal, 2016). Such dimensionality varied from
117 4,000 in pigs and chickens to 15,000 in Holsteins (Pocrnic et al., 2016c). Assuming limited
118 dimensionality, the inverse of \mathbf{G} (\mathbf{G}^{-1}) – as needed by REML – can be sparsely constructed using
119 the APY algorithm, with close to linear memory and computing requirements. Subsequently, the
120 inverses for over 2 million animals can be computed and stored (Tsuruta et al., 2021). However,
121 the inverse of \mathbf{H} also includes the inverse of a pedigree-based relationship matrix for genotyped
122 animals (Aguilar et al., 2010). Such a matrix can be dense with a long pedigree, but it is sparser
123 with a shorter pedigree. Thus, it could not be efficiently stored in large populations but had to be
124 accommodated indirectly (Strandén and Mäntysaari, 2014; Masuda et al., 2017).

125 The first purpose of this study was to find whether the costs of ssGREML can be reduced
126 using the APY algorithm with truncated pedigree and phenotypes. We hypothesize the truncation
127 could help to preserve the system's sparsity, given that APY \mathbf{G}^{-1} is sparser than the inverse of the

128 pedigree relationship matrices for deep pedigrees. The second purpose was to investigate to what
129 extent such truncation influences variance components and heritability estimates when different
130 numbers of core animals are used in APY.

131

132

Material and Methods

133 Animal care and use committee approval was not needed because data were simulated.

134

135 Data simulation

136 To evaluate the computational effectiveness of the proposed approach for estimating
137 variance components using genomic information, we simulated data using the QMSim software
138 (Sargolzaei et al., 2011). The simulator generated a historical population undergoing drift and
139 mutation and a recent population undergoing selection. The historical population consisted of
140 1,000 generations with a constant size of 50,000 individuals. Then, 800 more generations were
141 simulated where the number of individuals was reduced to 20,000, mimicking a bottleneck event.
142 The recent population (P1) consisted of 20 males and 15,000 females randomly sampled from
143 the last historical generation based on high phenotypic values. Individuals were mated along ten
144 generations producing a litter size of 1 with an equal probability of being male or female,
145 following a random mating design. Moreover, we considered a sire replacement rate of 0.50 and
146 a dam replacement rate of 0.20. Genomic information was available for 45,000 animals from
147 generations 8 through 10 (three youngest generations).

148 A total of 29 chromosomes of different lengths (ranging from 40 to 146 cM) were
149 simulated. Biallelic markers ($n = 52,890$) were evenly spaced along the chromosomes with equal

150 frequency in the first generation of the historical population. Potentially, 1,242 quantitative trait
151 loci (QTL) affected the trait and explained all the additive genetic variation; QTL allele effects
152 were sampled from a Gamma distribution with a shape parameter of 0.4. The mutation rate for
153 markers (recurrent mutation) and QTL was assumed to be equal to 2.5×10^{-5} per locus per
154 generation (Solberg et al., 2008).

155 The simulated trait had phenotypic variance and mean of 1.0, heritability and QTL
156 heritability of 0.30, and residual variance of 0.70. The simulated phenotypes were composed of:

$$157 \quad \mathbf{y} = \boldsymbol{\mu} + \mathbf{u} + \mathbf{e}$$

158 where \mathbf{y} is the vector of phenotypes, $\boldsymbol{\mu}$ is the vector of overall mean, \mathbf{u} is the vector of weighted
159 sum of QTL effects (i.e., additive genetic effect or animal effect), \mathbf{e} is the vector of residuals.

160 The standard error of estimates was small using 5 replicates during preliminary investigations of
161 this study. Because of that, the results are based on one replicate.

162

163 **Variance components**

164

165 Variance components were estimated using the average information (AI) REML
166 algorithm as implemented in the AIREMLF90 software (Misztal et al., 2002), which was
167 modified to incorporate the YAMS package (Masuda et al., 2014; Masuda et al., 2015). The
168 incorporation of YAMS was essential for this kind of task when using genomic information. The
169 package applies the supernodal method using multi-core optimized libraries (i.e., parallel
170 computing). The most computationally expensive part of the variance component estimation is
171 obtaining the inverse of the coefficient matrix used in traces. To that, efficient algorithms are
172 used to invert large and sparse matrices, which are based on three steps (i) ordering, (ii)

173 factorization (i.e., symbolic and numerical), and (iii) sparse inversion. Ordering is not
174 mandatory, but it saves a large amount of memory and time in the factorization step as it reduces
175 the *fill-in* effect (zero elements in the original matrix could become nonzero elements in the
176 factorized matrix). This effect can be minimized by ordering using appropriate techniques. In the
177 next step, the coefficient matrix (LHS of the mixed model equations) is factorized into two
178 triangular matrices by LU decomposition – L matrix. Finally, the Takahashi algorithm can be
179 used for inversion. The supernodal method is expected to provide faster inversions because they
180 find and process dense blocks in sparse matrices. Note that LHS inversion is only required to
181 estimate variance components or compute prediction error variance (PEV, obtained from
182 diagonal elements of an inverted LHS). If the objective is to solve the system of equations to
183 obtain breeding values, iterative methods as the preconditioned conjugate gradient (Lidauer et
184 al., 1999; Tsuruta et al., 2001) can be efficiently applied.

185 The model used to estimate variance components was based on the single-step method, in
186 which the inverse of the realized relationship matrix (\mathbf{H}^{-1}) is used in the mixed model equations
187 instead of \mathbf{A}^{-1} . Single-step genomic BLUP (ssGBLUP) is used for breeding value estimation,
188 whereas ssGREML is used for variance components estimation. The inversion of \mathbf{H} is computed
189 as follows (Aguilar et al., 2010):

190

$$191 \quad \mathbf{H}^{-1} = \mathbf{A}^{-1} + \begin{bmatrix} 0 & 0 \\ 0 & \mathbf{G}_{APY}^{-1} - \mathbf{A}_{22}^{-1} \end{bmatrix}$$

192

193 where \mathbf{A}^{-1} is the inverse of the pedigree relationship matrix, \mathbf{A}_{22}^{-1} is the inverse of the pedigree
194 relationship matrix for genotyped animals, computed by the algorithm described in Colleau
195 (2002). The genomic relationship matrix (\mathbf{G}) was computed as follows:

196

197

$$\mathbf{G} = \frac{\mathbf{Z}\mathbf{Z}'}{2 \sum p_j(1 - p_j)}$$

198

199 where \mathbf{Z} is the matrix of gene content centered by the current allele frequencies, and p_j is the
 200 allele frequency of SNP j . Inbreeding coefficients were considered when constructing the three
 201 relationship matrices. This provides a better equivalence between genomic and pedigree-based
 202 relationship matrices, leading to a more similar genetic base (Aguilar et al., 2020). The $\mathbf{G}_{\text{APY}}^{-1}$ is
 203 the inverse of the genomic relationship matrix obtained using the algorithm for proven and
 204 young (APY) (Misztal et al., 2014; Misztal, 2016). This algorithm considers that genotyped
 205 individuals are arbitrarily divided into core (c) and noncore (n). Breeding values for noncore
 206 (\mathbf{u}_n) can be described as a linear function of breeding values of core (\mathbf{u}_c):

207

208

$$\mathbf{u}_n = \mathbf{P}_n \mathbf{u}_c + \Phi_n$$

209

210 where $\mathbf{P}_n = \mathbf{Z}_n(\mathbf{Z}'_c \mathbf{Z}_c + \mathbf{I}\alpha)^{-1} \mathbf{Z}'_c$ is a matrix that relates breeding values of noncore and core,
 211 and Φ_n is the mendelian sampling term which has non-diagonal variance but can be
 212 approximated to diagonal. In cases where the number of core is large enough, breeding values of
 213 noncore depend only on breeding values of core (see Misztal (2016) for additional details). The
 214 inverse of \mathbf{G}_{APY} is constructed as following:

215

216

$$\mathbf{G}_{\text{APY}}^{-1} = \begin{bmatrix} \mathbf{I} - \mathbf{P}'_{cc} & -\mathbf{P}_{cn} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{M}_{cc}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{M}_{nn}^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{I} - \mathbf{P}_{cc} & \mathbf{0} \\ -\mathbf{P}_{nc} & \mathbf{I} \end{bmatrix}$$

217

218 If $\mathbf{G}_{cc}^{-1} = (\mathbf{I} - \mathbf{P}'_{cc})\mathbf{M}_{cc}^{-1}(\mathbf{I} - \mathbf{P}_{cc})$ is known, the complete inverse can be simplified to:

219

$$220 \quad \mathbf{G}_{APY}^{-1} = \begin{bmatrix} \mathbf{G}_{cc}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} + \begin{bmatrix} -\mathbf{P}_{cn} \\ \mathbf{I} \end{bmatrix} \mathbf{M}_{nn}^{-1} \begin{bmatrix} -\mathbf{P}_{nc} & \mathbf{I} \end{bmatrix}$$

221

222 where $\mathbf{P}_{cc} = \mathbf{G}_{cc}\mathbf{G}_{cc}^{-1}$, $\mathbf{M}_{cc(nn)} = \text{diag}\{g_{i,i} - p_{i,1:i-1}g'_{i,1:i-1}\}$ for individual i in the core

223 (noncore) group. Because \mathbf{G}_{APY}^{-1} is conditioned only on the genotypic information of core

224 animals, the matrix is sparser than the full \mathbf{G}^{-1} regularly used in ssGBLUP (Misztal, 2016). Note

225 that the covariance between two noncore individuals is null, but variances are stored in the

226 matrix.

227 The construction of the genomic matrix using APY in BLUPF90 software can be done in

228 two possible implementations. The first construction builds a single matrix for all core and

229 noncore. The second construction builds the genomic matrix in blocks and it aims to save

230 computing memory as it require less operations than single matrix (Masuda et al., 2016).

231 Currently, the single matrix construction is implemented for variance component estimation.

232

233 **Scenarios**

234

235 The scenarios below were built to evaluate the impact of the (1) size of the core group in

236 APY and the (2) influence of skipping zero elements from the LHS under different amounts of

237 pedigree and phenotypic data used in variance components estimation.

238

239 *Core group of different sizes*

240

241 Pocrnic et al. (2016a) evaluated the prediction accuracy using APY in simulation tests.
242 The authors suggested the greatest accuracy was found by selecting the number of core
243 individuals equal to the number of largest eigenvalues explaining 98% of \mathbf{G} (a number from now
244 on referred to as eigen98). This study tested core groups of different sizes to evaluate the impact
245 on variance components and heritability estimates. A total of four scenarios were tested by
246 allocating 1K (one thousand), 5K, 9K, and 14K randomly sampled out of 45,000 genotyped
247 individuals. For each of those scenarios, the largest variation explained was 72.03% (eigen70),
248 91.09% (eigen90), 95.70% (eigen95), and 98.07% (eigen98), respectively. For computational
249 reasons, the singular value decomposition of \mathbf{Z} was calculated instead of the eigenvalue
250 decomposition of \mathbf{G} .

251

252 *Evaluating the influence of pedigrees and phenotypes*

253

254 Using $\mathbf{G}_{\text{APY}}^{-1}$ helps to reduce computing time for genomic predictions because of its
255 sparsity (Fragomeni et al., 2015; Masuda et al., 2016); however, in the single-step approach, the
256 combined \mathbf{H}^{-1} contains also \mathbf{A}^{-1} and \mathbf{A}_{22}^{-1} , which are relatively dense. The APY method was
257 earlier applied to the construction of \mathbf{A}_{22}^{-1} without success (Breno Fragomeni, personal
258 communication). Although the sparsity of \mathbf{A}_{22}^{-1} may not be a requirement for genomic
259 predictions, it becomes essential for reducing computing time for variance components
260 estimation to follow the sparsity of $\mathbf{G}_{\text{APY}}^{-1}$. A reduction in the number of generations was
261 attempted to increase the sparsity in \mathbf{A}^{-1} and \mathbf{A}_{22}^{-1} . A total of seven different scenarios were
262 designed, differing on the number of pedigree generations used for variance components
263 estimation. Reduction in the generations of phenotypes was also used to follow pedigree

264 incompleteness and avoid bias. The scenarios were designed to mimic a real situation where the
265 actual founder population is usually unknown. Only three genotyped generations (45,000 most
266 recent animals) were kept in the genomic file for further analyses. Subsequent scenarios were
267 constructed by removing one generation of phenotypes and pedigree at a time, from the oldest to
268 the youngest animals.

269

270 *The influence of zero elements in the Mixed Model Equations (MME)*

271

272 Lastly, a scenario aimed to evaluate the impact of discarding zero elements from the LHS
273 of MME on computing performance and variance components estimation. For that, the *OPTION*
274 *skip_zero_in_dense_matrix* was used in AIREMLF90 (Misztal et al., 2014) to store only non-
275 zero elements of $\mathbf{G}_{APY}^{-1} - \mathbf{A}_{22}^{-1}$. When this option was used, the scenario was termed “Reduced”,
276 and otherwise “Full”.

277

278 RESULTS AND DISCUSSION

279 Previous studies have investigated the properties of APY, including its implementation
280 for large-scale genomic evaluations (Fragomeni et al., 2015; Lourenco et al., 2015; Masuda et
281 al., 2016) and its efficiency in real and simulated populations with different effective population
282 sizes (Pocrnic et al., 2016b; Pocrnic et al., 2016c). Bradford et al. (2017) studied the impact of
283 different core definitions, and Misztal et al. (2020) evaluated the GEBV fluctuation when
284 changing the core group in APY. Additionally, Vandenplas et al. (2018) investigated the impact
285 of using APY on GEBV estimation in crossbreeding schemes; Hidalgo et al. (2021) compared

286 the GEBV variation due to the inclusion of new data and changing the APY core animals.
287 Finally, Lourenco et al. (2018) studied the impact of using $\mathbf{G}_{\text{APY}}^{-1}$ instead of \mathbf{G}^{-1} on the estimation
288 of SNP effects. Our study evaluated the feasibility of using APY for variance components
289 estimation, the impact of removing generations of pedigree and phenotypic data on computing
290 time, and the influence of using a different number of core animals to construct the genomic
291 matrix. Variance components were estimated using AIREML modified to incorporate the YAMS
292 package for sparse matrix calculations (Masuda et al., 2014).

293

294 **Heritability estimates and computing performance**

295

296 Heritability, residual variance, and additive variance estimated using a different number
297 of generations in the pedigree and cores sizes in APY are shown in Figures 1-3. The standard
298 deviation of variance components and heritability across generations is shown in Table 1.
299 Because the simulation involved a certain level of selection, the expected heritability should
300 slightly deviate from the simulated value of 0.3. Therefore, the scenario with 10 generations of
301 data (i.e., full pedigree and full phenotypes) was used as a benchmark.

302 In general, the variance components and heritability estimates approached the simulated
303 values as the number of core approached eigen98. The scenario using 1K individuals (i.e.,
304 eigen70) in the core was the most sensitive to removing generations, suggesting that variance
305 components are highly impacted when the core group in APY represents the number of
306 eigenvalues explaining a smaller fraction of the total variation in \mathbf{G} . From a prediction accuracy
307 standpoint, a similar behavior was also observed in other studies (Pocrnic et al., 2016a; Pocrnic
308 et al., 2016c); however, the impact on variance components had not been investigated before.

309 Although pedigrees were more limited after removing a few generations of data, the combination
310 of pedigree and genomic information and the use of a proper core size controlled the bias in
311 variance components and heritability estimation. Small fluctuations on variance components and
312 heritability were observed when retaining only 4 to 6 generations of pedigree and phenotypes
313 with a core size equal to eigen98. In these scenarios, the difference in heritability was almost
314 nonexistent; this was also true when comparing Full and Reduced models.

315 The ratio σ_e^2/σ_a^2 is important when predicting breeding values using the mixed model
316 equations as it is the shrinkage factor for additive effects. The variability of the ratio under
317 different core sizes is shown in Figure 4. As the core size approached eigen98, the ratio became
318 closer to the simulated value of 2.33. Additionally, the ratio became less influenced by the
319 number of generations used to estimate the variance components as the core size approached
320 eigen98. Reliable variance components estimates (or at least their ratio and heritability) are of
321 great importance to ensure the accurate prediction of breeding values. The resulting breeding
322 values are not BLUP unless the true variances are known or are approaching the true parameters
323 (Kennedy, 1981).

324 The adoption of a core group that explains less than eigen98 affected the ability to
325 represent all the independent chromosome segments segregating in the population, traceback
326 gene frequencies, and consequently, accurately establish covariances between genotypic values.
327 In this study, we might have three different sources of changes for genetic variances. The first
328 source is related to the lack of relationships because generations were sequentially removed in
329 different scenarios. Unknown relationships (i.e., incorrect base population definition) affect the
330 estimation of Mendelian sampling variance in different intensities depending on the number of
331 known parents. If both parents are unknown, Mendelian sampling is equal to $0.5\sigma_a^2$, and if only

332 one parent is known, it equals $[0.75 - 0.25 \times f_p] \sigma_a^2$, where f_p is the inbreeding coefficient
333 (Henderson, 1976). Under mixed models, offspring breeding values are estimated as a function
334 of parent breeding values and Mendelian sampling. Thus, all individuals with unknown
335 relationships are treated as samples from the base population with average breeding value of 0
336 and common variance σ_a^2 .

337 The second source of change in genetic variance is the presence of selection over
338 generations, which affects the distribution of sire and dam breeding values. Unfortunately, it is
339 impossible to identify the contribution of each factor separately because this study was not
340 designed for that purpose. The third source of genetic variation, which is the aim of this study, is
341 the intentional use of a sparse representation of \mathbf{G}^{-1} , i.e., \mathbf{G}_{APY}^{-1} . In APY, it is intrinsically
342 assumed that the complete genome is divided into many independent chromosome segments
343 (ICS) containing non-redundant genomic information. The number of ICS is a statistical concept
344 that depends on the effective population size and the genome length (Stam, 1980). The
345 consequence of this assumption is that a small error in variance components estimation can be
346 observed by building the core group considering the dimensionality of \mathbf{G} as a function of the
347 number of eigenvalues explaining a certain proportion of variance. For example, if \mathbf{G}_{APY}^{-1} is built
348 based on the number of core animals equal to that of eigenvalues explaining 98% of the variance
349 in \mathbf{G} , the assumed error is 2% (Miształ et al., 2020). Results from the current study add a new
350 dimension to the factors driving the estimation of reliable variance components in the genomic
351 era. Thus, if the definition of the core group considers the genetic architecture of the population,
352 \mathbf{G} might contain all the genetic information necessary to estimate reliable variance components
353 (Junqueira et al., 2017; Junqueira et al., 2020). In addition to the factors evaluated in this study,

354 Cesarani et al. (2019) have found that the selection design and genotyping structure can
355 influence bias in estimating variance components.

356

357 **Computing resources**

358

359 Nowadays, much effort has been placed on developing faster and computationally
360 feasible methods for a virtually unlimited number of genotyped individuals. Using large-scale
361 datasets becomes more problematic when the objective is to estimate variance components. This
362 is because most algorithms require several rounds of inversion of the LHS of MME before the
363 convergence is reached. During computations, factorization and inversion are the most
364 demanding steps in the REML estimation. The possibility to combine APY to compute a sparse
365 representation of \mathbf{G}^{-1} , data reduction, and YAMS (i.e., dense blocks operation) (Masuda et al.,
366 2014; Masuda et al., 2015) seems computationally beneficial. In this study, we evaluated the
367 factors impacting the timing required for computational operations. Figure 5 shows the average
368 computing time, relative to total (i.e., in percentage), required for ordering, factorization
369 (symbolic and numerical), and sparse inversion with data reduction (pedigree and phenotypes).
370 The most time-consuming operation was the inversion, which took more than 50% of the total
371 time. This was expected because matrix inversion has a cubic computing cost. Next, numerical
372 factorization consumed nearly 30% of the total computing time, whereas ordering and symbolic
373 factorization took approximately 9% and 7.5%, respectively. Skipping zero elements in the
374 MME did not improve the computing time of any of the inverse operations.

375 A detailed description of the computing time required by each step after data removal is
376 in Figure 6. The descriptive statistics of computing time savings across generations is shown in

377 Table 2. Ordering showed the most prominent timing decrease due to data removal, followed by
378 symbolic factorization among the four steps. On average, a 7% decrease in the computing time
379 for ordering was observed by removing each generation of data. During MME computations,
380 ordering and symbolic factorization are not mandatory. These operations are mainly
381 implemented to reduce computing time for numerical factorization and inversion. As more
382 genotypes and/or pedigree records are included in the model, the time required for numerical
383 factorization and sparse inversion increases. Using a simulated dataset with \mathbf{G}_{APY}^{-1} and YAMS,
384 we observed an opposite behavior where shorter pedigree sometimes caused an increase in
385 computing time for the numerical factorization and sparse inversion operations. In these
386 operations, there were no gains in computing performance due to data removal, as shown by the
387 regression slope, which was close to 0 (Table 2). The greatest savings were around 10% when
388 using six generations of pedigree and phenotypic data. It is known that numerical factorization
389 and sparse inversion are the most demanding operations in REML computations. The fact that
390 the required time for these operations was not reduced can be explained by the creation of
391 nonzero elements not present in the coefficient matrix before the numerical factorization is done.
392 Those elements are known as “fill-in elements.”

393 Consequently, extra calculations are needed, obviously increasing the amount of time to
394 complete the sparse inversion. There are several efforts in developing faster algorithms focused
395 on typical nonzero structures in sparse matrices. The sparse matrix algorithm in YAMS uses
396 supernodal techniques (i.e., common nonzero pattern between adjacent columns) to speed-up
397 computations. Computing time might be significantly improved compared to other sparse matrix
398 packages (e.g., FSPAK) because the memory hierarchy is more effectively exploited in dense

399 operations, and multiple columns within a submatrix are simultaneously updated (Masuda et al.,
400 2014).

401 **Conclusions**

402 The algorithm for proven and young (APY) can be successfully applied to create the
403 inverse of the genomic relationship matrix used in single-step genomic restricted maximum
404 likelihood for estimating variance components. To ensure reliable variance component
405 estimation, it is important to use a core size that corresponds to the number of largest eigenvalues
406 explaining around 98% of total variation in \mathbf{G} . When APY is used, pedigrees can be truncated to
407 increase the sparsity of \mathbf{H} and slightly reduce computing time for ordering and symbolic
408 factorization, with no impact on the estimates. A reduction in computing time for numerical
409 factorization and sparse inversion is unlike because of fill-in elements. The savings in
410 computing time for estimating variance components is far less than the expected efficiency that
411 APY has shown compared to the use of regular \mathbf{G}^{-1} for breeding values estimation. This
412 inefficiency is because the block implementation of APY is still not possible for variance
413 components estimation.

414

415 **Conflict of interest statement**

416 The authors declare that they do not have any conflict of interest.

417

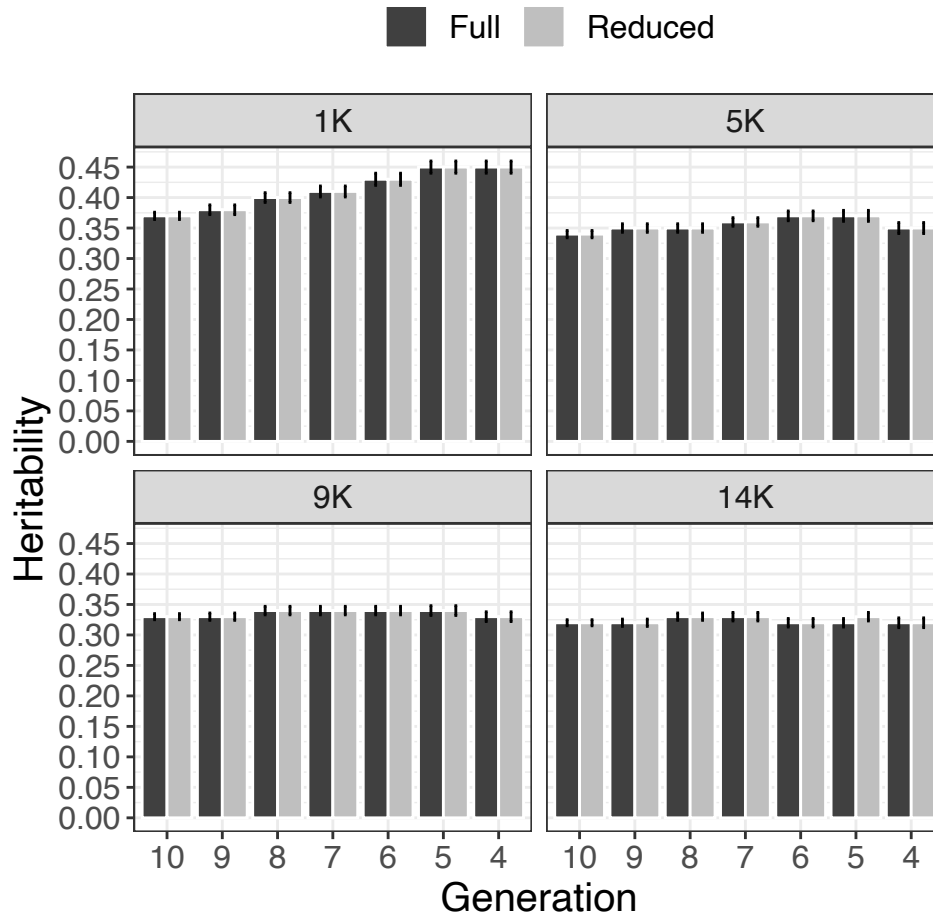
418 **Literature Cited**

- 419 Aguilar, I., E. N. Fernandez, A. Blasco, O. Ravagnolo, and A. Legarra. 2020. Effects of ignoring
420 inbreeding in model-based accuracy for BLUP and SSGBLUP. *Journal of Animal*
421 *Breeding and Genetics* 137(4):356-364.
- 422 Aguilar, I., I. Misztal, D. L. Johnson, A. Legarra, S. Tsuruta, and T. J. Lawlor. 2010. Hot topic:
423 A unified approach to utilize phenotypic, full pedigree, and genomic information for
424 genetic evaluation of Holstein final score. *Journal of Dairy Science* 93(2):743-752. doi:
425 10.3168/jds.2009-2730
- 426 Bradford, H., I. Pocnić, B. Fragomeni, D. Lourenco, and I. Misztal. 2017. Selection of core
427 animals in the algorithm for proven and young using a simulation model. *Journal of*
428 *Animal Breeding and Genetics* 134(6):545-552.
- 429 Cesarani, A., G. Gaspa, F. Correddu, M. Cellesi, C. Dimauro, and N. Macciotta. 2019. Genomic
430 selection of milk fatty acid composition in Sarda dairy sheep: Effect of different
431 phenotypes and relationship matrices on heritability and breeding value accuracy. *Journal*
432 *of Dairy Science* 102(4):3189-3203.
- 433 Colleau, J. J. 2002. An indirect approach to the extensive calculation of relationship coefficients.
434 *Genetics Selection Evolution* 34(4):409.
- 435 Fragomeni, B. O., D. A. L. Lourenco, S. Tsuruta, Y. Masuda, I. Aguilar, A. Legarra, T. J.
436 Lawlor, and I. Misztal. 2015. Hot topic: Use of genomic recursions in single-step
437 genomic best linear unbiased predictor (BLUP) with a large number of genotypes.
438 *Journal of Dairy Science* 98(6):4090-4094. doi: 10.3168/jds.2014-9125
- 439 Henderson, C. R. 1975. Best linear unbiased estimation and prediction under a selection model.
440 *Biometrics*:423-447.
- 441 Henderson, C. R. 1976. A simple method for computing the inverse of a numerator relationship
442 matrix used in prediction of breeding values. *Biometrics*:69-83.
- 443 Hidalgo, J., D. Lourenco, S. Tsuruta, Y. Masuda, S. Miller, M. Bermann, A. L. Garcia, and I.
444 Misztal. 2021. Changes in genomic predictions when new information is added. *Journal*
445 *of Animal Science* 99(2):skab004.
- 446 Junqueira, V. S., F. F. Cardoso, M. M. Oliveira, B. P. Sollero, F. F. Silva, and P. S. Lopes. 2017.
447 Use of molecular markers to improve relationship information in the genetic evaluation
448 of beef cattle tick resistance under pedigree-based models. *Journal of Animal Breeding*
449 *and Genetics* 134(1):14-26.
- 450 Junqueira, V. S., P. S. Lopes, D. Lourenco, F. F. e Silva, and F. F. Cardoso. 2020. Applying the
451 metafounders approach for genomic evaluation in a multibreed beef cattle population.
452 *Frontiers in Genetics* 11
- 453 Kennedy, B. 1981. Variance component estimation and prediction of breeding values. *Canadian*
454 *Journal of Genetics and Cytology* 23(4):565-578.

- 455 Lidauer, M., I. Strandén, E. A. Mäntysaari, J. Pösö, and A. Kettunen. 1999. Solving large test-
456 day models by iteration on data and preconditioned conjugate gradient. *J. Dairy Sci.*
457 82(12):2788-2796.
- 458 Lourenco, D. A., B. O. Fragomeni, S. Tsuruta, I. Aguilar, B. Zumbach, R. J. Hawken, A.
459 Legarra, and I. Misztal. 2015. Accuracy of estimated breeding values with genomic
460 information on males, females, or both: an example on broiler chicken. *Genetics*
461 *Selection Evolution* 47(1):56.
- 462 Lourenco, D. A. L., A. Legarra, S. Tsuruta, D. Moser, S. Miller, and I. Misztal. 2018. Tuning
463 indirect predictions based on SNP effects from single-step GBLUP. *Interbull Bulletin*
464 (53)
- 465 Masuda, Y., I. Aguilar, S. Tsuruta, and I. Misztal. 2015. Technical note: Acceleration of sparse
466 operations for average-information REML analyses with supernodal methods and sparse-
467 storage refinements. *Journal of Animal Science* 93(10):4670-4674.
- 468 Masuda, Y., T. Baba, and M. Suzuki. 2014. Application of supernodal sparse factorization and
469 inversion to the estimation of (co) variance components by residual maximum likelihood.
470 *Journal of Animal Breeding and Genetics* 131(3):227-236.
- 471 Masuda, Y., I. Misztal, A. Legarra, S. Tsuruta, D. A. L. Lourenco, B. O. Fragomeni, and I.
472 Aguilar. 2017. Avoiding the direct inversion of the numerator relationship matrix for
473 genotyped animals in single-step genomic best linear unbiased prediction solved with the
474 preconditioned conjugate gradient. *Journal of Animal Science* 95(1):49-52.
- 475 Masuda, Y., I. Misztal, S. Tsuruta, A. Legarra, I. Aguilar, D. A. L. Lourenco, B. O. Fragomeni,
476 and T. J. Lawlor. 2016. Implementation of genomic recursions in single-step genomic
477 best linear unbiased predictor for US Holsteins with a large number of genotyped
478 animals. *Journal of Dairy Science* 99(3):1968-1974.
- 479 Meyer, K. 1997. An average information restricted maximum likelihood algorithm for estimating
480 reduced rank genetic covariance matrices or covariance functions for animal models with
481 equal design matrices. *Genetics Selection Evolution* 29(2):97.
- 482 Misztal, I. 2016. Inexpensive computation of the inverse of the genomic relationship matrix in
483 populations with small effective population size. *Genetics* 202(2):401-409.
- 484 Misztal, I., H. L. Bradford, D. A. L. Lourenco, S. Tsuruta, Y. Masuda, A. Legarra, and T. J.
485 Lawlor. 2017. Studies on inflation of GEBV in single-step GBLUP for type. *Interbull*
486 *Bulletin* (51):38-42.
- 487 Misztal, I., A. Legarra, and I. Aguilar. 2014. Using recursion to compute the inverse of the
488 genomic relationship matrix. *Journal of Dairy Science* 97(6):3943-3952.
- 489 Misztal, I., S. Tsuruta, I. Pocrnic, and D. Lourenco. 2020. Core-dependent changes in genomic
490 predictions using the algorithm for proven and young in single-step genomic best linear
491 unbiased prediction. *Journal of Animal Science* 98(12):skaa374.

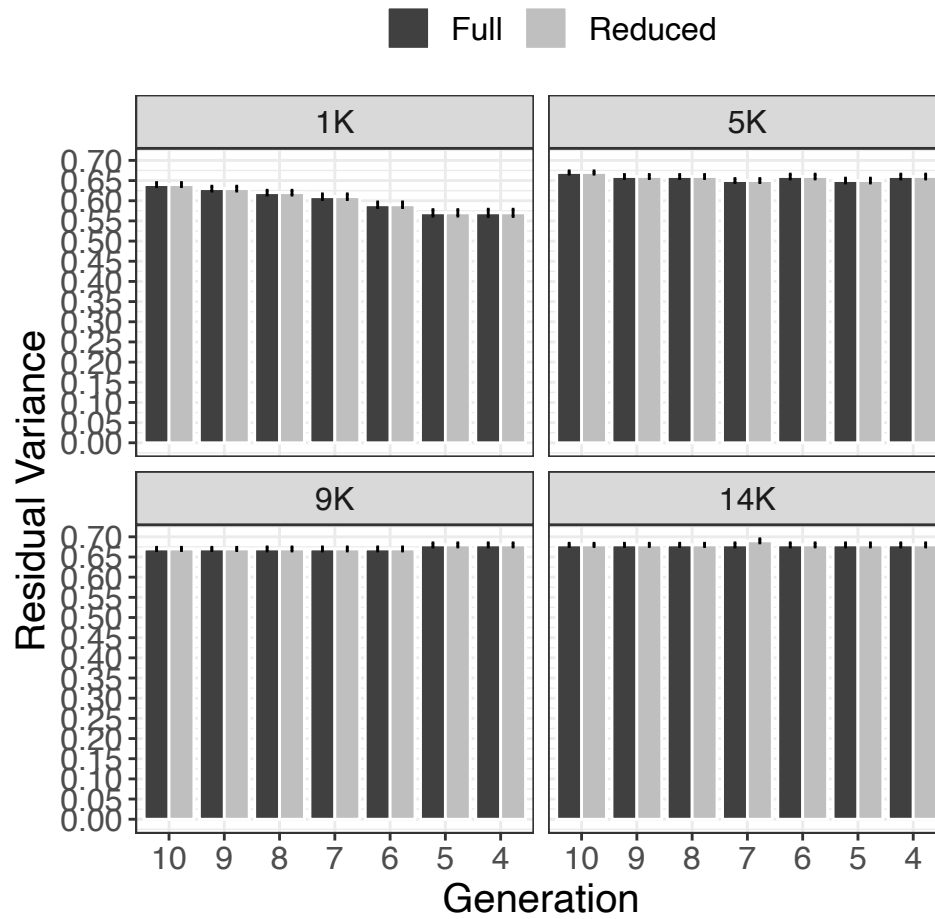
- 492 Misztal, I., S. Tsuruta, T. Strabel, B. Auvray, T. Druet, and D. H. Lee. 2002. BLUPF90 and
493 related programs. In: Proceedings of the 7th World Congress on Genetics Applied to
494 Livestock Production
- 495 Patry, C., and V. Ducrocq. 2011. Evidence of biases in genetic evaluations due to genomic
496 preselection in dairy cattle. *Journal of Dairy Science* 94(2):1011-1020.
- 497 Patterson, H. D., and R. Thompson. 1971. Recovery of inter-block information when block sizes
498 are unequal. *Biometrika* 58(3):545-554.
- 499 Pocrnic, I., D. A. Lourenco, Y. Masuda, A. Legarra, and I. Misztal. 2016a. The dimensionality of
500 genomic information and its effect on genomic prediction. *Genetics* 203(1):573-581.
- 501 Pocrnic, I., D. A. L. Lourenco, Y. Masuda, A. Legarra, and I. Misztal. 2016b. The
502 dimensionality of genomic information and its effect on genomic prediction. *Genetics*
503 203(1):573-581.
- 504 Pocrnic, I., D. A. L. Lourenco, Y. Masuda, and I. Misztal. 2016c. Dimensionality of genomic
505 information and performance of the Algorithm for Proven and Young for different
506 livestock species. *Genetics Selection Evolution* 48(1):82.
- 507 Sargolzaei, M., J. Chesnais, and F. Schenkel. 2011. FImpute-An efficient imputation algorithm
508 for dairy cattle populations. *Journal of Dairy Science* 94(1):421.
- 509 Solberg, T., A. Sonesson, J. Woolliams, and T. Meuwissen. 2008. Genomic selection using
510 different marker types and densities. *Journal of Animal Science* 86(10):2447-2454.
- 511 Stam, P. 1980. The distribution of the fraction of the genome identical by descent in finite
512 random mating populations. *Genetics Research* 35(2):131-155.
- 513 Strandén, I., and E. A. Mäntysaari. 2014. Comparison of some equivalent equations to solve
514 single-step GBLUP. In: Proceedings of the 10th World Congress on genetics applied to
515 Livestock production. Vancouver. p 22.
- 516 Tsuruta, S., T. Lawlor, D. Lourenco, and I. Misztal. 2021. Bias in genomic predictions by mating
517 practices for linear type traits in a large-scale genomic evaluation. *Journal of Dairy*
518 *Science* 104(1):662-677.
- 519 Tsuruta, S., I. Misztal, and I. Strandén. 2001. Use of the preconditioned conjugate gradient
520 algorithm as a generic solver for mixed-model equations in animal breeding applications.
521 *J. Anim. Sci.* 79(5):1166-1172.
- 522 Vandenplas, J., M. P. Calus, and J. Ten Napel. 2018. Sparse single-step genomic BLUP in
523 crossbreeding schemes. *Journal of Animal Science* 96(6):2060-2073.
- 524 VanRaden, P. M. 2008. Efficient methods to compute genomic predictions. *Journal of Dairy*
525 *Science* 91(11):4414-4423.
526

527
528



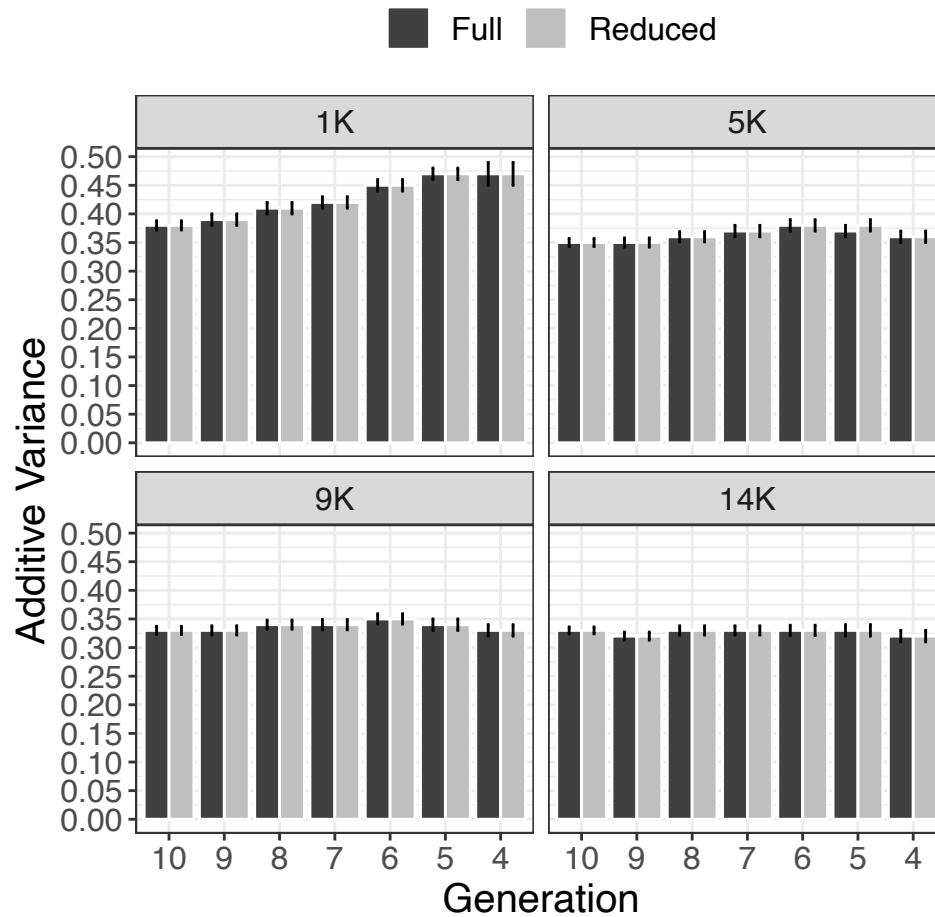
529
530
531
532
533

Figure 1. Heritability calculated along one replicate of simulation considering different number of generations with pedigree and phenotypic data under different number of core animal in APY. Two scenarios were considered, where zeros were stored (Full) or not (Reduced). Error bars represent the standard error of prediction under REML.



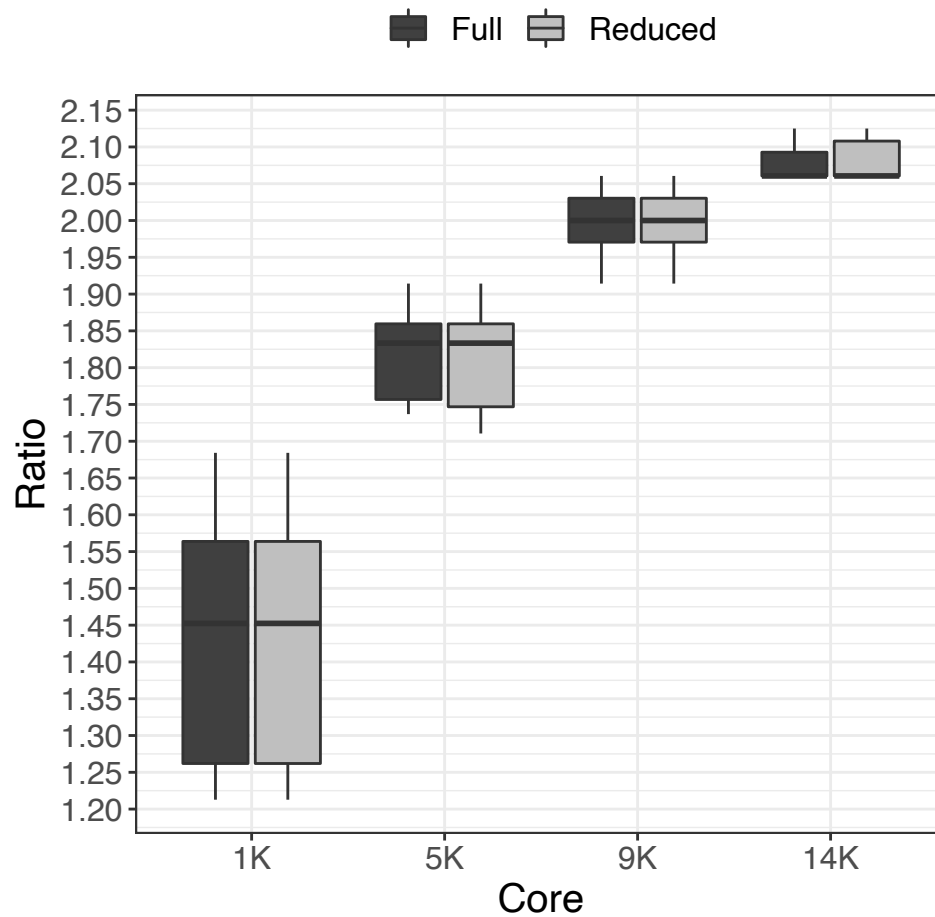
534
535
536
537
538
539
540
541

Figure 2. Residual variance calculated along one replicate of simulation considering different number of generations with pedigree and phenotypic data under different number of core animal for APY calculation. Two scenarios were considered, where zeros were stored (Full) or not (Reduced). Error bars represent the standard error of prediction under REML.



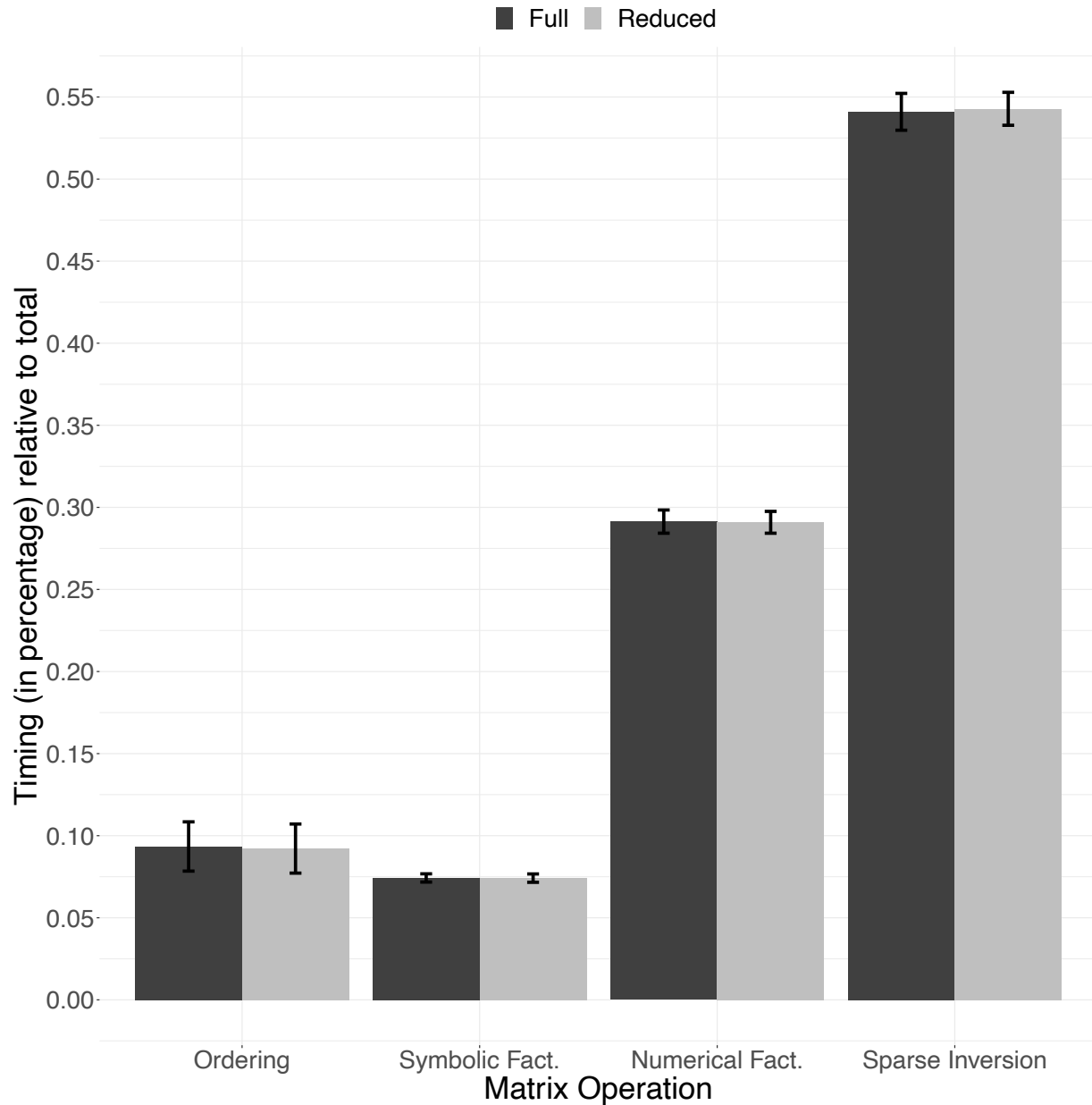
542
543
544
545
546
547
548
549

Figure 3. Additive variance calculated along one replicate of simulation considering different number of generations with pedigree and phenotypic data under different number of core animal in APY. Two scenarios were considered, where zeros were stored (Full) or not (Reduced). Error bars represent the standard error of prediction under REML.

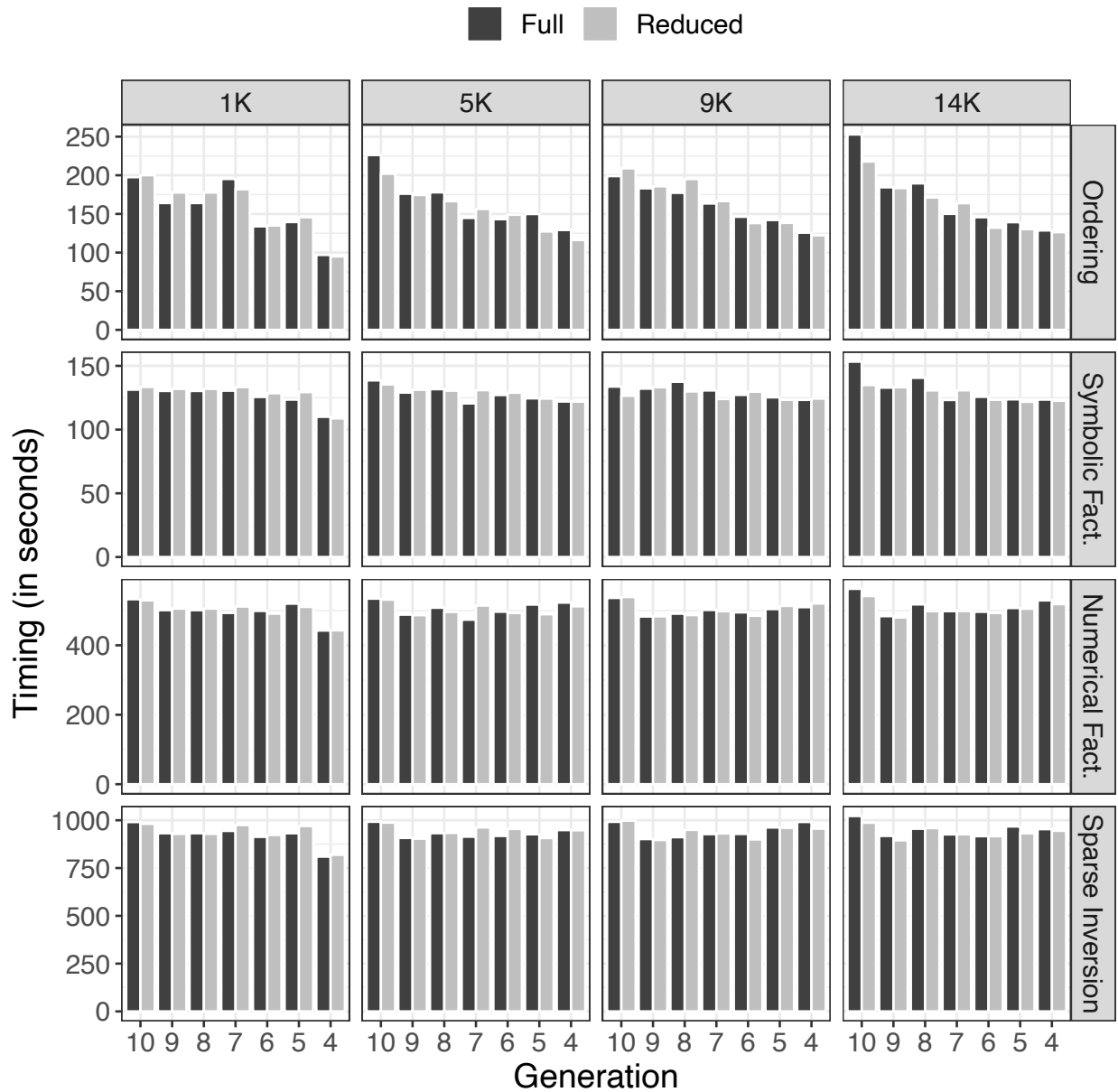


550
551
552
553
554
555
556

Figure 4. Distribution of the ratio (σ_e^2 / σ_a^2) over different number of generations with pedigree and phenotypic data using different sizes for the core group in APY. Two scenarios were considered, where zeros were stored (Full) or not (Reduced). Error bars represents the standard error of prediction under REML.



557
558 **Figure 5.** Average timing in percentage (ratio between total timing) relative to each
559 operation used in the process of matrix inversion. The average timing and error bars
560 (standard deviation) were calculated across scenarios using different number of
561 generations in the pedigree and phenotypic and core sizes. The x-axis represents the
562 steps required to invert matrices: finding the ordering, symbolic factorization (Symbolic
563 Fact., setting up the data structure), numerical factorization (Numerical Fact.), and
564 sparse inversion. Two scenarios were considered, where zeros were stored (Full) or not
565 (Reduced).



566
 567 **Figure 6.** Timing (in seconds) relative to each operation to invert matrices using
 568 different number of generations in the pedigree and phenotypes under different number
 569 of core animals for the computation of $APY \mathbf{G}^{-1}$. Matrix inversion steps: finding the
 570 ordering (Ordering), symbolic factorization (Symbolic Fact.), numerical factorization
 571 (Numerical Fact.), and sparse inversion. Two scenarios were considered, where zeros
 572 were stored (Full) or not (Reduced).
 573
 574

575 **Table 1.** Standard deviation of variance components and heritability calculated across
576 generations using a complete (Full) mixed model equations (MME), and a reduced
577 MME after skipping zero elements (Reduced).
578

Parameter ¹	Core ²	Scenario	
		Full	Reduced
σ_a^2	1K	0.037	0.037
	5K	0.011	0.013
	9K	0.008	0.008
	14K	0.005	0.005
σ_e^2	1K	0.028	0.028
	5K	0.007	0.007
	9K	0.005	0.005
	14K	0.000	0.004
h^2	1K	0.032	0.032
	5K	0.011	0.011
	9K	0.005	0.005
	14K	0.005	0.005

579 ¹ σ_a^2 : additive variance, σ_e^2 : residual variance, h^2 : heritability

580 ² Number of individuals included as core to build the inverse of genomic matrix using the
581 algorithm of proven and young (APY)

582

Table 2. Descriptive statistics of computing time savings for the matrix operations and the slope of a regression of computing time on generations after removing pedigree and phenotypic data. The benchmark is the model using full pedigree and phenotypic data. The comparison is based on using core group of different sizes in algorithm for proven and young (APY), and based on a full mixed model equations (Full) and a reduced mixed model equations after skipping zero elements (Reduced).

Core Size	Matrix Operation	Full						Reduced					
		Min (%)	Mean (%)	Max (%)	SD (%) ¹	Slope ²	³	Min (%)	Mean (%)	Max (%)	SD (%)	Slope	
1K	Ordering	1.16	24.58	50.94	16.98	-0.07	**	9.10	23.95	52.47	16.99	-0.08	**
	Symbolic Factorization	0.61	4.77	16.22	6.04	-0.02	**	0.08	4.57	18.44	6.92	-0.02	*
	Numerical Factorization	2.38	7.47	16.92	4.93	-0.02	ns	3.36	6.57	16.32	4.97	-0.02	*
	Sparse Inversion	4.67	8.08	18.25	5.08	-0.02	**	0.59	5.80	16.47	5.71	-0.01	ns
5K	Ordering	21.35	32.13	42.88	8.60	-0.06	**	13.61	26.58	42.48	11.19	-0.07	**
	Symbolic Factorization	4.98	9.24	13.06	3.08	-0.02	**	3.10	5.50	9.97	2.89	-0.01	**
	Numerical Factorization	2.13	6.21	11.34	3.47	-0.00	ns	3.22	6.17	8.41	2.22	-0.00	**
	Sparse Inversion	4.39	6.81	8.52	1.48	-0.00	ns	2.52	5.33	8.51	2.51	-0.00	ns
9K	Ordering	7.94	21.40	36.80	11.13	-0.06	**	6.65	24.55	41.48	13.99	-0.07	ns
	Symbolic Factorization	2.76	6.39	10.33	2.97	-0.01	**	2.48	5.19	7.43	2.20	-0.01	ns
	Numerical Factorization	4.96	7.26	9.98	1.82	-0.00	ns	3.45	7.63	10.33	2.93	-0.00	ns
	Sparse Inversion	0.07	5.52	9.08	3.38	-0.00	ns	3.77	6.59	10.14	2.82	-0.00	ns
14K	Ordering	25.04	38.19	49.13	9.85	-0.07	**	15.79	30.57	41.97	11.27	-0.07	**
	Symbolic Factorization	8.28	16.33	19.63	4.61	-0.03	**	1.26	5.79	9.71	3.72	-0.02	**
	Numerical Factorization	5.92	10.15	13.99	2.89	-0.00	ns	4.32	7.91	11.39	2.35	-0.00	ns
	Sparse Inversion	5.34	8.09	10.32	2.14	-0.01	ns	2.85	5.88	9.32	2.25	-0.00	ns

¹ Standard deviation

² Slope of a regression of computing time on generations

³ Slope statistical significance * P<0.05, **P<0.10, ^{ns}not significant

bioRxiv preprint doi: <https://doi.org/10.1101/2022.01.19.476983>; this version posted January 21, 2022. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY-NC-ND 4.0 International license](#).

1

2