

1

2 SINGLE STEP, A GENERAL APPROACH FOR GENOMIC SELECTION

3

4 Andres Legarra^{a*}, Ole F Christensen^b, Ignacio Aguilar^c and Ignacy Misztal^d

5

6 ^aINRA, UMR1388 GenPhySE, BP52627, 31326 Castanet Tolosan, France

7 andres.legarra@toulouse.inra.fr

8 ^bCenter for Quantitative Genetics and Genomics, Department of Molecular Biology and

9 Genetics, Aarhus University, Blichers Alle 20, P.O. BOX 50, DK-8830 Tjele, Denmark

10 OleF.Christensen@agrsci.dk

11 ^cInstituto Nacional de Investigación Agropecuaria, Las Brujas 90200, Uruguay

12 iaguilar@inia.org.uy

13 ^dDepartment of Animal and Dairy Science, University of Georgia, Athens 30602-2771, USA

14 ignacy@uga.edu

15 *Corresponding author: andres.legarra@toulouse.inra.fr

16 Phone:+33561285182 Fax: +33561285353

17

18

19 ABSTRACT

20 Genomic evaluation methods assume that the reference population is genotyped and
21 phenotyped. This is most often false and the generation of pseudo-phenotypes is uncertain and
22 inaccurate. However, markers obey transmission rules and therefore the covariances of
23 marker genotypes across individuals can be modelled using pedigree relationships. Based on
24 this, an extension of the genomic relationship matrix can be constructed in which genomic
25 relationships are propagated to all individuals, resulting in a combined relationship matrix,
26 which can be used in a BLUP procedure called the Single Step Genomic BLUP. This
27 procedure provides so far the most comprehensive option for genomic evaluation. Several
28 extensions, options and details are described: compatibility of genomic and pedigree
29 relationships, Bayesian regressions, multiple trait models, computational aspects, etc. Many
30 details scattered through a series of papers are put together into this paper.

31

32

33 KEYWORDS

34 Genetic evaluation, genomic evaluation, marker genotypes, BLUP, relationship

35

37 **1. INTRODUCTION: BRIEF EXCURSION INTO METHODS FOR GENOMIC**

38 **EVALUATION**

39 *1.1 Marker information*

40 Genetic progress by selection and mating is based on prediction of the ability of the parents to
41 breed the most efficient descendants. This process of prediction is called genetic evaluation or
42 prediction. Genetic evaluation in plants and livestock has, for the last century, been based on
43 the use of phenotypes at the traits of interest, together with pedigree. In most cases, these
44 evaluations ignore the physical base of heredity, i.e., DNA, and use a simplified conceptual
45 representation of the transmission of genetic information from parents to offspring; namely,
46 each parent passes on average half its genetic constitution, associated with an unknown
47 sampling known as Mendelian sampling. Recent technical developments allow stepping
48 further into biology and peering at the genome in the form of single nucleotide
49 polymorphisms, known as SNP markers. These markers depict, in an incomplete manner, the
50 differences between DNA inherited by two individuals. They can be used in multiple ways; in
51 this section we will present very briefly how they are typically used in genetic evaluation (or
52 prediction or estimation of breeding values: EBV hereinafter) in a parametric framework.
53 Most genomic evaluations follow the principle of estimating the *conditional expectation* of
54 the breeding value in view of all information, which has optimal properties if the assumptions
55 of the model hold (e.g., Fernando and Gianola 1986). This (parametric) paradigm has been
56 extremely fruitful over the last decades, allowing for the development of BLUP, REML,
57 Bayesian estimators and giving a coherent framework to solve many applied problems in
58 animal breeding (e.g., Gianola and Fernando, 1986).

59 The notion of prediction or estimation of random effects is absent in many statistical
60 textbooks (but check, for instance, Casella and Berger (1990)). However, it has been treated
61 as early as Smith (1936) with key references e.g. in Cochran (1951), Henderson (1973) or
62 Fernando and Gianola (1986). Based on those authors, the “correct” model of prediction
63 consists in writing down the statistical association between phenotypes and breeding values,
64 then derive the EBVs from the conditional distribution of breeding values given the
65 phenotypes. |

66

67 *1.2 Bayesian regression*

68 Typically, in genomic predictions, the phenotypes of a population are considered as a function
69 of the breeding values, and the breeding value of individuals, \mathbf{u} (or part of it) is decomposed
70 into a sum of marker effects \mathbf{a} (e.g., Meuwissen et al., 2001; VanRaden, 2008). These marker
71 effects are summed according to the genotype of the individual, coded as (0,1,2) for the
72 (AA, Aa, aa) genotypes. In matrix notation $\mathbf{u} = \mathbf{M}\mathbf{a}$. It follows that one way of estimating
73 breeding values is to estimate marker effects and then use $\hat{\mathbf{u}} = \mathbf{M}\hat{\mathbf{a}}$. In order to estimate
74 marker effects, one needs to assume a prior distribution for them. The process of estimation of
75 marker effects using the statistical model for phenotypes $p(\mathbf{y}|\mathbf{a})$ and the prior for markers
76 $p(\mathbf{a})$ is often called *Bayesian Regression on markers*. A difficult decision is the choice of the
77 prior for markers. An extensive literature in the subject shows higher accuracy, for some traits
78 and populations, of using “heavy-tailed” a priori distributions (e.g., VanRaden et al., 2009).

79

80 *1.3 RR-BLUP or GBLUP*

81 If multivariate normality is assumed for the effect of markers, interesting things happen in the
82 algebraic developments. The first one is that the **Bayesian Regression** becomes what is called
83 RR-BLUP (or SNP-BLUP). The second is the existence of closed forms for the **RR-BLUP**
84 estimators of marker effects, in the form of Henderson's Mixed Model Equations; these
85 estimators greatly simplify computations and can be easily extended, e.g. for multiple trait
86 situations. The third is the existence of a so-called equivalent model, in which breeding values
87 (and not marker effects) are directly computed by **Henderson's Mixed Model Equations** using
88 a covariance matrix $Var(\mathbf{u}) = \mathbf{ZD}_a\mathbf{Z}'$ (VanRaden, 2008), where $\mathbf{Z} = \mathbf{M} - 2\mathbf{P}$ and \mathbf{P} contains
89 p_k , the allelic frequencies of markers. **This is most often called GBLUP.** In the most common
90 case it is assumed that $Var(\mathbf{a}) = \mathbf{D}_a = \mathbf{I}\sigma_u^2/2\sum p_k q_k$, where σ_u^2 is the genetic variance, so
91 that that $Var(\mathbf{u}) = \sigma_u^2\mathbf{G}$, where $\mathbf{G} = \mathbf{Z}\mathbf{Z}'/2\sum p_k q_k$. The matrix \mathbf{G} is called the *genomic*
92 *relationship matrix* and will frequently be referred to later. Properties of \mathbf{G} for populations in
93 Hardy-Weinberg equilibrium are an average diagonal of 1 and an average off-diagonal of 0.
94 Genomic evaluation using \mathbf{G} (**GBLUP**) gives the same estimated breeding values as a marker-
95 based RR-BLUP and has the additional advantage of fitting very well into ancient
96 developments (e.g., for multiple trait) and current software. An interesting feature of the
97 genomic relationship matrix is that it can be seen as an "improved" estimator of relationships
98 based on markers instead of pedigrees (VanRaden, 2008; Hayes *et al.*, 2009), and is closely
99 related to estimators of relationships based on markers used in conservation genetics (Ritland,
100 1996; Toro *et al.*, 2011).

101

102 **2. THE PROBLEM OF MISSING GENOTYPES AND THE USE OF PSEUDO-DATA**

103 Genotyping an individual is an expensive process that also requires the availability of a
104 biological sample. Therefore, in most populations either the most recent or the most

105 representative animals (e.g., sires in dairy cattle) have been genotyped. Some individuals are
106 genotyped with low-density chips that genotype only some markers. From these, genotypes at
107 all markers can be efficiently imputed (e.g., VanRaden et al., 2013) and we will consider
108 these individuals as genotyped. A non-genotyped individual is one for which *there is no*
109 *genotype at any loci*. Therefore, the methods for genomic prediction described above cannot
110 be applied directly, as there is often not phenotype for the individual genotyped and viceversa;
111 this is particularly true for sex-limited traits (milk yield, fertility, prolificacy). **Although a sire**
112 **model could be used, this ignores selection on the female side, and does not yield females'**
113 **EBVs**. Therefore, animal breeders have used pseudo-data or *pseudo-phenotypes*. A pseudo-
114 phenotype is a projection of the phenotypes of individuals close to the genotyped one. In dairy
115 cattle and sheep, pseudo-phenotypes typically used are corrected daughter performances
116 (daughter yield deviations, VanRaden and Wiggans, 1991), whereas in other species de-
117 regressed proofs are often used, with a variety of *ad hoc* adjustments (Garrick et al., 2009;
118 Ricard et al., 2013).

119 This process is therefore clumsy and we call it *multiple step*. A regular **genetic** evaluation
120 based on pedigree is run first, and its results are used to create pseudo-performances. Then, a
121 genomic evaluation model is used. This results in losses of information, inaccuracies and
122 biases, whose importance depends on the species and data set. There are several possible
123 problems:

- 124 1. The information of a close relative is ignored in the genomic prediction, for instance
125 the dam of a bull if this dam has phenotype but not genotype.
- 126 2. The information of a close relative is ignored in the creation of pseudo-phenotypes, for
127 instance a non-genotyped parent. This is serious if the progeny of the genotyped
128 individual is scarce and therefore parental phenotypes are informative (see Ricard et
129 al. (2013) for a discussion in a horse application).

- 130 3. Unless estimates of environmental effects are perfect, covariances among pseudo-
131 phenotypes are not correctly modelled. For instance, the yield deviations of two
132 unrelated cows in the same herd will be correlated (e.g., if the herd effect is
133 underestimated both will be biased upwards). This is ignored in the genomic model,
134 which acts as if pseudo-phenotypes were perfectly clean of environmental errors.
- 135 4. Many key parameters are difficult to obtain. One of them is precisions of pseudo-
136 phenotypes, which are in most cases rough approximations.
- 137 5. There is no feedback. An improved estimation of the breeding value of the genotyped
138 animal should go into the regular pedigree-based genetic evaluation and improve its
139 global accuracy.
- 140 6. When genomic selection is applied, animals are selected as parents based on their
141 known genotype. The implication is that when phenotypes are obtained from a scheme
142 that has used genomic selection, evaluation based on pedigree becomes biased and is
143 no longer appropriate (Patry and Ducrocq, 2011). Hence, current approaches for
144 constructing pseudo-phenotypes will also become inappropriate due to problems of
145 bias.
- 146 7. The process is extremely difficult to generalize. For instance, the multiple-trait
147 generalization of pseudo-phenotypes is basically non-existent, and the pseudo-
148 phenotypes for maternal traits result in much less accurate multiple step predictions
149 (Lourenco et al., 2013).

150 Some of these defaults can be palliated. VanRaden et al. (2009) used a selection index to *a*
151 *posteriori* add information from non-genotyped dams to bull genomic evaluations. The
152 procedures of creation of pseudo-phenotypes can be refined over and over, and in dairy cattle
153 they result in very accurate predictions, as accurate as Single Step (Aguilar *et al.*, 2010). In
154 other species the adequacy of multiple step procedures varies more. However, the existence of

155 these problems calls for a unified procedure for prediction of genetic value. This paper will
156 describe such a procedure: the *Single Step*.

157

158 **3. DEVELOPMENT OF THE SINGLE STEP METHOD FOR GENOMIC** 159 **EVALUATION**

160 Legarra et al. (2009) and Christensen and Lund (2010) developed in parallel the basic theory
161 for the Single Step. They started from two somehow different points of view that turned out to
162 result in the same formulation, and we will present both developments, starting with the latter
163 one.

164

165 *3.1 The Single Step as “imputing” missing genotypes*

166 To some extent, missing genotypes can be deduced from existing genotypes, for instance a
167 dam mated to a sire AA producing an offspring Aa is necessarily carrier of one allele a . In
168 statistical theory, a way to deal with missing information is to augment the model with this
169 missing information (*e.g.*, Tanner and Wong, 1987). This missing information needs to be
170 inferred from the other data, and its joint distribution needs to be considered. This means that
171 a “best guess” of missing information in view of observed data, as suggested by Hickey et al.
172 (2012), who imputed genotypes for the complete nongenotyped population, is not correct
173 enough. Even if one considers the uncertainty of individual “guesses” the across-individual
174 uncertainty is extremely difficult to ascertain or deal with.

175 An example may clarify this point. Assume a very long complex pedigree and the final
176 generation genotyped for one locus, with allelic frequency $p = frequency(a)$. Due to only
177 having one generation with genotypes and to the long and complex pedigree, best guesses of

178 genotypes in the base animals will be nearly identical and equal to $2p$, for all individuals.
179 Therefore, using “best guess” of genotype without taking uncertainty into account, all base
180 population individuals will be treated by the genomic evaluation as identical, which will force
181 them to have the same estimated breeding value, which is paradoxical. For each individual the
182 uncertainty can be assessed by noting that the distribution of genotypes in this case is
183 approximately AA (with probability q^2), Aa (with probability $2pq$) and aa (with probability
184 p^2), but the joint distribution of genotypes for individuals in the base population is much
185 more difficult to characterize. In principle, incorporation of uncertainty can be done by
186 sampling all possible genotypic configurations of all individuals, e.g. by a Gibbs sampling
187 procedure (e.g. Abraham et al., 2007) but this is computationally infeasible for data of the size
188 used in practical genetic evaluations.

189

190 Christensen and Lund (2010), considered the problem as follows. Their objective was to
191 create an extension of the genomic relationship matrix to nongenotyped animals. Following
192 an idea of Gengler et al. (2007), they treated the genotypes as quantitative traits. This makes
193 sense because genotypes are quantitative (0/1/2) and follow Mendelian transmissions.
194 Therefore the covariance of the genotypes z of two individuals i and j is described by their
195 relationship, i.e. $Cov(z_i, z_j) = A_{ij}2pq$ (e.g., Cockerham, 1969). This is less informative than
196 considering the genotype as a union of two discrete entities following Mendelian rules (e.g.,
197 sometimes we can exactly deduce a genotype from close relatives) but makes the problem
198 analytically tractable for all cases.

199

200 Christensen and Lund (2010) started by inferring the genomic relationship matrix for all
201 animals using inferred (imputed) genotypes for nongenotyped animals; these can simply be

202 obtained as $\hat{\mathbf{Z}}_1 = \mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{Z}_2$, where 1 and 2 stand for nongenotyped and genotyped animals,
 203 respectively. This provides the “best guess” of genotypes. However, the missing data theory
 204 requires the joint distribution of these “guessed” genotypes. Assuming that multivariate
 205 normality holds for genotypes (this is an approximation, but very good when many genotypes
 206 are considered), the “best guess” is $E(\mathbf{Z}_1|\mathbf{Z}_2) = \hat{\mathbf{Z}}_1$, and the conditional variance expressing
 207 the uncertainty about the “guess” is $Var(\hat{\mathbf{Z}}_1|\mathbf{Z}_2) = (\mathbf{A}_{11} - \mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{A}_{21})\mathbf{V}$ where \mathbf{V} contains
 208 $2p_kq_k$ (where $q_k = 1 - p_k$) in the diagonal. These two results can be combined to obtain the
 209 desired augmented genomic relationships. For instance, for the nongenotyped animals,

$$Var(\mathbf{u}_1) = \sigma_u^2 \left(\frac{\hat{\mathbf{Z}}_1\hat{\mathbf{Z}}_1'}{2\sum p_kq_k} + \mathbf{A}_{11} - \mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{A}_{21} \right),$$

210 which equals

$$Var(\mathbf{u}_1) = \sigma_u^2(\mathbf{A}_{11} - \mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{A}_{21} + \mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{G}\mathbf{A}_{22}^{-1}\mathbf{A}_{21})$$

211 Finally, the augmented covariance matrix is

$$Var \begin{pmatrix} \mathbf{u}_1 \\ \mathbf{u}_2 \end{pmatrix} = \sigma_u^2 \mathbf{H},$$

212 where

$$\mathbf{H} = \begin{pmatrix} \mathbf{A}_{11} - \mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{A}_{21} + \mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{G}\mathbf{A}_{22}^{-1}\mathbf{A}_{21} & \mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{G} \\ \mathbf{G}\mathbf{A}_{22}^{-1}\mathbf{A}_{21} & \mathbf{G} \end{pmatrix},$$

213 is the augmented genomic relationship matrix with inverse

$$\mathbf{H}^{-1} = \mathbf{A}^{-1} + \begin{pmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{G}^{-1} - \mathbf{A}_{22}^{-1} \end{pmatrix}$$

214 assuming that \mathbf{G} is invertible (this will be dealt with later). Therefore, by using an algebraic

215 data augmentation of missing genotypes, Christensen and Lund (2010) derived a simple

216 expression for an augmented genomic relationship matrix and its inverse, without the need to
217 explicitly augment, or “guess”, all genotypes for all non-genotyped animals.

218

219 *3.2 The Single Step as Bayesian updating of the relationship matrix*

220 Legarra et al. (2009) arrived to the same expressions that of Christensen and Lund (2010) in a
221 different manner. They also considered how to construct an extended relationship matrix.

222 However, instead of dealing with individual markers, they dealt with overall breeding values

223 that can be written as $\mathbf{u}_2 = \mathbf{Z}_2 \mathbf{a}$. They reasoned as follows. Prior to observation of markers,

224 the joint distribution of breeding values is multivariate normal

$$p \begin{pmatrix} \mathbf{u}_1 \\ \mathbf{u}_2 \end{pmatrix} = N(\mathbf{0}, \sigma_u^2 \mathbf{A})$$

225 with covariance matrix

$$\text{Var} \begin{pmatrix} \mathbf{u}_1 \\ \mathbf{u}_2 \end{pmatrix} = \sigma_u^2 \mathbf{A} = \sigma_u^2 \begin{pmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{pmatrix}$$

226 After observing the markers, this covariance matrix will change. The joint distribution above

227 can be split into the product of a marginal and a conditional density; i.e. $p(\mathbf{u}_1, \mathbf{u}_2) =$

228 $p(\mathbf{u}_1 | \mathbf{u}_2) p(\mathbf{u}_2)$, where

$$229 \quad p(\mathbf{u}_1 | \mathbf{u}_2) = N(\mathbf{A}_{12} \mathbf{A}_{22}^{-1} \mathbf{u}_2, \sigma_u^2 (\mathbf{A}_{11} - \mathbf{A}_{12} \mathbf{A}_{22}^{-1} \mathbf{A}_{21})).$$

230 In other terms, $\mathbf{u}_1 = \mathbf{A}_{12} \mathbf{A}_{22}^{-1} \mathbf{u}_2 + \boldsymbol{\epsilon}$, where $\boldsymbol{\epsilon}$ and \mathbf{u}_2 are independent, and $\text{Var}(\boldsymbol{\epsilon}) =$

$$231 \quad \sigma_u^2 (\mathbf{A}_{11} - \mathbf{A}_{12} \mathbf{A}_{22}^{-1} \mathbf{A}_{21}).$$

232 As discussed before, in presence of marker genotypes the genomic relationship matrix can be

233 considered as fully informative about relationships of individuals, without the need to resort

234 to pedigree or knowledge of previous, or future, nongenotyped individuals. Therefore, *after*
 235 observing the marker genotypes

$$p(\mathbf{u}_2 | \text{markers}) = N(\mathbf{0}, \sigma_u^2 \mathbf{G}).$$

236 Marker genotypes influence the relationships among nongenotyped individuals and
 237 relationships between nongenotyped and genotyped individuals indirectly. Assuming that
 238 these relationships are only influenced by marker genotypes through the genomic
 239 relationships among genotyped individuals, and assuming that the statistical distribution is
 240 determined by these relationships, one can write that

$$p(\mathbf{u}_1 | \mathbf{u}_2, \text{markers}) = p(\mathbf{u}_1 | \mathbf{u}_2)$$

241 Therefore, the joint distribution of breeding values *after* observing the markers is:

$$p(\mathbf{u}_1, \mathbf{u}_2 | \text{markers}) = p(\mathbf{u}_1 | \mathbf{u}_2) p(\mathbf{u}_2 | \text{markers})$$

242 From these results, expressions for the covariance of breeding values are immediate. For
 243 instance, $Var(\mathbf{u}_1) = \sigma_u^2 (\mathbf{A}_{12} \mathbf{A}_{22}^{-1} \mathbf{G} \mathbf{A}_{22}^{-1} \mathbf{A}_{21} + \mathbf{A}_{11} - \mathbf{A}_{12} \mathbf{A}_{22}^{-1} \mathbf{A}_{21})$, where the part involving \mathbf{G}
 244 is the variability associated to the conditional mean of breeding values of nongenotyped
 245 individuals given the genotyped ones; and the second part is the variability beyond this
 246 conditional mean. Finally, the result

$$Var \begin{pmatrix} \mathbf{u}_1 \\ \mathbf{u}_2 \end{pmatrix} = \sigma_u^2 \mathbf{H} = \sigma_u^2 \begin{pmatrix} \mathbf{A}_{11} - \mathbf{A}_{12} \mathbf{A}_{22}^{-1} \mathbf{A}_{21} + \mathbf{A}_{12} \mathbf{A}_{22}^{-1} \mathbf{G} \mathbf{A}_{22}^{-1} \mathbf{A}_{21} & \mathbf{A}_{12} \mathbf{A}_{22}^{-1} \mathbf{G} \\ \mathbf{G} \mathbf{A}_{22}^{-1} \mathbf{A}_{21} & \mathbf{G} \end{pmatrix}$$

247 is obtained, in full agreement with Christensen and Lund (2010). The reason for this
 248 agreement is that in both cases a central assumption is that the influence of marker genotypes
 249 on nongenotyped individuals is via relationships determined by the numerator relationship
 250 matrix \mathbf{A} .

251

252 3.3 Genetic properties of the extended relationship matrix

253 Matrix \mathbf{H} above can be seen as a modification of regular pedigree relationships to
254 accommodate genomic relationships. For instance, two seemingly unrelated individuals will
255 appear as related in \mathbf{H} if their descendants are related in \mathbf{G} . Accordingly, two descendants of
256 individuals that are related in \mathbf{G} will be related in \mathbf{H} , even if the pedigree disagrees. Indeed, it
257 has been suggested (Sun et al., 2013) to use \mathbf{H} in mating programs to avoid inbreeding.

258 Contrary to common intuition from BLUP or GBLUP, genotyped animals without phenotype
259 or descendants *cannot* be eliminated from matrix \mathbf{H} . The reason is that (unless both parents
260 are genotyped) these animals potentially modify pedigree relationship across other animals,
261 possibly notably their parents. For instance imagine two half-sibs, offspring of one sire mated
262 to two nongenotyped, unrelated cows. If these two half sibs are virtually identical, \mathbf{H} will
263 include this information and the cows will be made related (even identical) in \mathbf{H} .

264

265 3.4 Single Step Genomic BLUP

266 Because the Single Step relationship matrix provides an explicit and rather sparse inverse of
267 the extended relationship matrix \mathbf{H} , its application to genomic evaluation is immediate. A full
268 specification of the Single Step Genomic BLUP assumes the following model:

$$\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{W}\mathbf{u} + \mathbf{e}$$

$$\text{Var}(\mathbf{u}) = \mathbf{H}\sigma_u^2; \text{Var}(\mathbf{e}) = \mathbf{I}\sigma_e^2$$

269 with \mathbf{H} and its inverse as shown above. The logic of BLUP (Henderson, 1973 and many other
270 publications) holds and the only change is to use \mathbf{H} instead of the numerator relationship

271 matrix. Genomic predictions estimating simultaneously all breeding values and using all
272 available information are, for the single trait case, the solutions to the mixed model equations
273 (e.g., Aguilar et al., 2010; Christensen and Lund, 2010):

$$\begin{pmatrix} \mathbf{X}'\mathbf{X} & \mathbf{X}'\mathbf{W} \\ \mathbf{W}'\mathbf{X} & \mathbf{W}'\mathbf{W} + \mathbf{H}^{-1}\lambda \end{pmatrix} \begin{pmatrix} \hat{\mathbf{b}} \\ \hat{\mathbf{u}} \end{pmatrix} = \begin{pmatrix} \mathbf{X}'\mathbf{y} \\ \mathbf{W}'\mathbf{y} \end{pmatrix}$$

274 where $\lambda = \sigma_e^2/\sigma_u^2$.

275

276 Note that any formulation using relationship matrix \mathbf{A} can use \mathbf{H} instead, and therefore there
277 is also Single Step REML and Single Step Gibbs, for instance in Legarra et al. (2011a) and
278 Forni et al. (2011).

279

280 **4. EXTENSIONS AND REFINEMENTS OF THE SINGLE STEP**

281 As said above, any model that has been fit as BLUP can be fit as Single Step. We will
282 describe a few of these extensions that are of interest.

283 *4.1 Pseudo-Single Step.*

284 Also called “blending” (e.g. Su et al., 2012a), this has been used to include all males of a
285 population with pseudo-phenotypes, where some are genotyped and some are not. This is a
286 compromise between using all information (which might be complex) and ignoring pseudo-
287 phenotypes of non-genotyped males, for instance sires of genotyped males. Accuracy
288 increases, but less than with true Single Step (Baloche et al., 2014).

289 *4.2 Multiple trait*

290 Extension to deal with multiple traits is immediate. The mixed model equations are, in the
 291 usual notation:

$$\begin{pmatrix} \mathbf{X}'\mathbf{R}^{-1}\mathbf{X} & \mathbf{X}'\mathbf{R}^{-1}\mathbf{W} \\ \mathbf{W}'\mathbf{R}^{-1}\mathbf{X} & \mathbf{W}'\mathbf{R}^{-1}\mathbf{W} + \mathbf{H}^{-1} \otimes \mathbf{G}_0 \end{pmatrix} \begin{pmatrix} \hat{\mathbf{b}} \\ \hat{\mathbf{u}} \end{pmatrix} = \begin{pmatrix} \mathbf{X}'\mathbf{R}^{-1}\mathbf{y} \\ \mathbf{W}'\mathbf{R}^{-1}\mathbf{y} \end{pmatrix}$$

292 where $\mathbf{R} = \mathbf{I} \otimes \mathbf{R}_0$, \mathbf{R}_0 is the matrix of residual covariances across traits and \mathbf{G}_0 is the matrix
 293 of genetic covariances across traits. Extension to random regressions or maternal effect
 294 models is very similar.

295

296 4.3 Marker effect estimates

297 The GBLUP and other models based on genomic relationship matrices such as the Single Step
 298 do not directly provide estimates of marker effects. These are of interest in order to spot
 299 locations of major genes (or QTL) and also in order to provide a less computationally
 300 demanding evaluation of new born animals that are genotyped but do not have phenotypes.
 301 The marker effects can be deduced from estimated breeding values of the genotyped
 302 individuals. Consider the joint distribution of breeding values \mathbf{u} and marker effects \mathbf{a}
 303 (Henderson, 1973; Strandén and Garrick, 2009):

$$Var \begin{pmatrix} \mathbf{u}_2 \\ \mathbf{a} \end{pmatrix} = \begin{pmatrix} \mathbf{Z}_2 \mathbf{D}_a \mathbf{Z}_2' & \mathbf{Z}_2 \mathbf{D}_a \\ \mathbf{D}_a \mathbf{Z}_2' & \mathbf{D}_a \end{pmatrix}$$

304 where, usually, $\mathbf{D}_a = \mathbf{I} \sigma_u^2 / 2 \sum p_i q_i$ (this assumption will be relaxed later). Assuming
 305 multivariate normality, $\hat{\mathbf{u}}_2 | \hat{\mathbf{a}} = \mathbf{Z}_2 \hat{\mathbf{a}}$ (the breeding value is the sum of marker effects) and

306 $\hat{\mathbf{a}} | \hat{\mathbf{u}}_2 = \mathbf{D}_a \mathbf{Z}_2' (\mathbf{Z}_2 \mathbf{D}_a \mathbf{Z}_2')^{-1} \hat{\mathbf{u}}_2 = \mathbf{D}_a \mathbf{Z}_2' \mathbf{G}^{-1} \sigma_u^{-2} \hat{\mathbf{u}}_2$ where (as discussed in previous sections)

307 $\mathbf{Z}_2 \mathbf{D}_a \mathbf{Z}_2' = \mathbf{G} \sigma_u^2$, so that marker effects can be deduced by backsolving using the genomic

308 relationship matrix and markers' incidence matrix. This result has been used, e.g., by Wang et
309 al. (2012), and it will appear later in this paper.

310

311 4.4 Extra polygenic effect

312 It has been often argued that markers do not capture all genetic variation. This can be shown
313 by estimating variance assigned to markers and pedigree (e.g. Legarra et al., 2008) or because
314 some genomic evaluation procedures give better cross-validation results when an extra
315 polygenic term based exclusively on pedigree relationships is added (e.g. Su et al., 2012b).
316 The GBLUP (VanRaden, 2008) and the derivations in the Single Step can accommodate this
317 very easily (Aguilar et al., 2010; Christensen and Lund, 2010). Let us decompose the breeding
318 values of genotyped individuals in a part due to markers and a residual part due to pedigree,
319 $\mathbf{u}_2 = \mathbf{u}_{m,2} + \mathbf{u}_{p,2}$ with respective variances $\sigma_u^2 = \sigma_{u,m}^2 + \sigma_{u,p}^2$. It follows that $Var(\mathbf{u}_2) =$
320 $(\alpha\mathbf{G} + (1 - \alpha)\mathbf{A}_{22})\sigma_u^2$ where $\alpha = \sigma_{u,m}^2/\sigma_u^2$. Therefore, the simplest way is to create a
321 modified genomic relationship matrix \mathbf{G}_w (\mathbf{G} in Aguilar et al., 2010; \mathbf{G}_w in VanRaden, 2008
322 and Christensen and Lund, 2010) as $\mathbf{G}_w = \alpha\mathbf{G} + (1 - \alpha)\mathbf{A}_{22}$ and to plug this relationship
323 matrix in all the expressions before. This has the additional advantage of making \mathbf{G}_w
324 invertible, which is not guaranteed for \mathbf{G} . Equivalently, one can fit *two* random effects, one
325 \mathbf{u}_m with covariance matrix $\mathbf{H}\sigma_{u,m}^2$ and another \mathbf{u}_p with covariance matrix $\sigma_{u,p}^2$.

326

327 4.5 Compatibility of genomic and pedigree relationships

328 This is a key issue in genomic evaluation that has received small attention beyond Single Step
329 developers even though, as shown by Vitezica et al. (2011), it also affects multiple step
330 methods. The derivations above of Single Step mixed model equations include terms such as

331 $\mathbf{G} - \mathbf{A}_{22}$ and $\mathbf{G}^{-1} - \mathbf{A}_{22}^{-1}$. This suggests that \mathbf{G} and \mathbf{A}_{22} , the genomic and pedigree
332 relationship matrices, need to be compatible. It has been long known (e.g., Ritland 1996) that
333 relationships estimated from markers need to use allelic frequencies at the base populations;
334 otherwise a severe bias in the estimated relationships is observed (VanRaden 2008; Toro et
335 al., 2011). However, typically base population frequencies are unknown because pedigree
336 recording started before biological sampling of individuals. The two derivations of the Single
337 Step assume, either implicitly or explicitly, that the base frequencies are known. In the
338 derivation of Christensen and Lund (2010) the allele frequencies enter explicitly. In the
339 derivation of Legarra et al. (2009) the hypothesis is that the expected breeding value of the
340 genotyped population is 0. This hypothesis will be wrong if either there has been selection or
341 drift, which is commonly the case; the average breeding value will change, and the genetic
342 variance will be reduced. These problems were soon observed by analysis of real life data sets
343 (Chen et al., 2011b; Forni et al., 2011; Christensen et al., 2012) and verified by simulation
344 (Vitezica et al., 2011).

345

346 Several proposals exist so far to make pedigree and genomic relationships compatible. The
347 three first proposals “tune” matrix \mathbf{G} to make it compatible with \mathbf{A}_{22} , in the form $\mathbf{G}^* = a +$
348 $b\mathbf{G}$, where a can be understood as an “overall” relationship and b as a change in scale (or
349 genetic variance). VanRaden (2008) suggested a regression of observed on expected
350 relationships, minimizing the residuals of $a + b\mathbf{G} = \mathbf{A}_{22} + \mathbf{E}$. This reflects the fact that over
351 conceptual repetitions of our population (same pedigree but different meiosis and genotypes)
352 $E(\mathbf{G}) = \mathbf{A}_{22}$ if \mathbf{G} is the realized relationship and \mathbf{A}_{22} is the expected relationship (VanRaden,
353 2008; Hayes et al., 2009). This idea was generalized to several breed origins by Harris and
354 Johnson (2010). The distribution of \mathbf{E} is not homoscedastic (Hill and Weir, 2011; Garcia-
355 Cortes et al., 2013) and this precluded scholars from trying this approach because it would be

356 sensible to extreme values (Christensen et al., 2012), e.g., if many far relatives are included,
357 for which the deviations in \mathbf{E} can be very large. A second approach is to model the
358 distribution of the mean of genotyped individuals, i.e., to assume a unknown mean μ for
359 genotyped individuals: $p(\mathbf{u}_2) = N(\mathbf{1}\mu, \mathbf{G})$. This is a random variable: the effect of selection
360 or drift on the trait will vary from one conceptual repetition to another. One can equally write
361 $p(\mathbf{u}_2) = N(\mathbf{0}, \mathbf{G} + \mathbf{1}\mathbf{1}'Var(\mu))$ with μ integrated out. An unbiased method forces the
362 distribution of average values of breeding values ($\bar{\mathbf{u}}_2$) to be identical and therefore, the
363 adjustment uses $\mathbf{G}^* = a + b\mathbf{G}$ with $b = \mathbf{1}$ and $a = \bar{\mathbf{A}}_{22} - \bar{\mathbf{G}}$ where the bar implies average
364 across values of \mathbf{G} and \mathbf{A} . Although this models corrects the change due to genetic trend, it
365 does not consider the fact that there is a reduction in genetic variance from the base
366 population to the genotyped individuals considered in \mathbf{A}_{22} but not in \mathbf{G} . This problem has
367 been tackled twice. The first manner is to consider genotyped individuals as a subpopulation
368 of all individuals in the population and to use Wright's fixation index theory, which allows
369 putting relationships in any scale (Cockerham, 1969, 1973). Translated to our context (Powell
370 et al., 2010) this implies $a = \bar{\mathbf{A}}_{22} - \bar{\mathbf{G}}$ and $b = 1 - a/2$ (Vitezica et al., 2011). The value of a
371 can be understood as an overall within-population relationship within the genotyped
372 individuals, with respect to an older population whose genotypes are not observed. This
373 overall relationship cannot be estimated by \mathbf{G} for lack of base allele frequencies. The value of
374 $a/2$ can be understood as the "extra" decrease in genetic variance in a random mating
375 population of average relationship $\bar{\mathbf{A}}_{22}$. Christensen et al. (2012) remarked that the hypothesis
376 of random mating population is not likely for the group of genotyped animals, since they
377 would born in different years and some being descendants of others, and suggested to infer a
378 and b jointly based on the drift of the mean of the population (as in Vitezica et al., 2011) and
379 based on the expected genetic variance, which is encapsulated in the average inbreeding
380 observed in \mathbf{G} and \mathbf{A}_{22} . More formally, the empirical variance of breeding values: $S_{u_2}^2 =$

381 $\mathbf{u}'_2 \mathbf{u}_2 / n - (\bar{\mathbf{u}}_2)^2$ has an expectation $\left(\frac{\text{tr}(\mathbf{A}_{22})}{n} - \bar{\mathbf{A}}_{22}\right) \sigma_u^2$ or $\left(\frac{\text{tr}(\mathbf{G}^*)}{n} - \bar{\mathbf{G}}^*\right) \sigma_u^2$ where n is the
382 number of individuals. Forcing unbiasedness implies that a and b should be determined from
383 the system of two equations: $\frac{\text{tr}(\mathbf{G})}{n} b + a = \frac{\text{tr}(\mathbf{A}_{22})}{n}$ and $+b\bar{\mathbf{G}} = \bar{\mathbf{A}}_{22}$. In random mating
384 populations in Hardy-Weinberg equilibrium (for instance in large populations of dairy cattle
385 and sheep, where Hardy-Weinberg equilibrium approximately holds), it turns out that
386 $b = 1 - a/2$ as in Vitezica et al. (2011). If restricting the group of animals for which
387 compatibility is required to those that are born in a certain generation, the assumption of
388 random mating among those genotyped animals is not unreasonable to assume in many
389 livestock species. All these corrections utilize some estimate of the allelic frequencies to
390 construct \mathbf{G} , and using observed allele frequencies (either based on all genotyped animals, or
391 based on a subset born in a certain generation) is usually done.

392 Finally, Christensen (2012) suggested the opposite point of view, to “tune” \mathbf{A}_{22} to \mathbf{G} instead
393 of the opposite. Pedigrees are arbitrary and depend on the start of pedigree, whereas
394 genotypes at the markers are absolute. Allele frequencies, though, change all the time. He
395 modelled the likelihood of markers given the pedigree as a quantitative trait and then
396 integrated over the uncertain allele frequencies. This amounts to fix allele frequencies at 0.5
397 and introduce two extra parameters, γ and s . The γ parameter can be understood as the overall
398 relationship across the base population such that current genotypes are more likely, and
399 integrates the fact that the assumption of unrelatedness at the base population is false in view
400 of genomic results (two animals who share alleles at markers are related even if the pedigree
401 is not informative). More precisely, he devised a new pedigree relationship matrix, $\mathbf{A}(\gamma)$
402 whose founders have a relationship matrix $\mathbf{A}_{bas} = \gamma + \mathbf{I}(1 - \gamma/2)$. Parameter s , used in
403 $\mathbf{G} = \mathbf{Z}\mathbf{Z}'/s$ can be understood as the counterpart of $2\sum p_k q_k$ (heterozygosity of the markers) in
404 the base generation. Both parameters can be deduced from maximum likelihood. This model

405 is the only one which introduces all the complexities of pedigrees (former ones are based on
406 average relationships) but it has not been tested with real data so far (Christensen, 2012).

407

408 *4.6 Computational algorithms*

409 The use and development of the Single Step has been possible through the use of several state
410 of the art algorithms. Construction and inversion of matrix \mathbf{G} are **cubic processes**, and are
411 much optimized by the use of efficient algorithms and parallel computations (Aguilar et al.,
412 2011). Construction of matrix \mathbf{A}_{22} has been possible, for very large pedigrees, by the
413 algorithm of Colleau (2002) which uses Henderson's decomposition of $\mathbf{A} = \mathbf{TDT}'$ to devise a
414 "solving" that allows easy multiplication of $\mathbf{w} = \mathbf{Av}$ and computation of \mathbf{A}_{22} in **quadratic**
415 time (Aguilar et al., 2011).

416 Further, the use of the solver known as preconditioned conjugated gradients (PCG) allows an
417 easy programming to solve the Single Step mixed model equations. PCG proceeds by
418 repeated multiplications **(LHS)sol** where **sol** is the vector of unknowns. In practice, this
419 product is split into a part

$$\begin{pmatrix} \mathbf{X}'\mathbf{X} & \mathbf{X}'\mathbf{W} \\ \mathbf{W}'\mathbf{X} & \mathbf{W}'\mathbf{W} + \mathbf{A}^{-1}\lambda \end{pmatrix} \begin{pmatrix} \hat{\mathbf{b}} \\ \hat{\mathbf{u}} \end{pmatrix}$$

420 for which very efficient algorithms already exist (e.g. Strandén and Lidauer, 1999) and a part

$$(\mathbf{G}^{-1} - \mathbf{A}_{22}^{-1})\lambda \hat{\mathbf{u}}_2$$

421 which can be done very efficiently, in particular using parallelization.

422 In addition, some implementations of the Single Step have used unsymmetric equations to
423 avoid inversion of \mathbf{G} (Misztal et al., 2009; Aguilar et al., 2013), with solution by the Bi-

424 Conjugate Gradient Stabilized algorithm. Legarra and Ducrocq (2012) reviewed and
425 suggested implementations of the Single Step with view towards very large data sets such as
426 in dairy cattle. Problems of these data sets are twofold. First, current evaluations use very
427 sophisticated software, first for regular BLUP (e.g., random regressions), and later for
428 genomic evaluations (e.g., Bayesian regressions). Second, the large size of the data sets,
429 which may preclude inversion (and even construction) of \mathbf{G} . They suggested two main
430 alternatives: a non-symmetric system of equations with non-inverted \mathbf{A}_{22} and \mathbf{G} , and an
431 iterative procedure similar to the multiple step but in which results from genomic evaluations
432 would be reintroduced in the regular BLUP evaluation, and results from regular BLUP would
433 be “data” for the genomic evaluations. The non-symmetric system shows slow convergence
434 on large data sets (Aguilar et al., 2013), whereas the iterative method is still untested on large
435 data sets. This is still an active field of research.

436

437 *4.7 Bayesian regressions in the Single Step*

438 Bayesian or non-linear regressions with non-normal priors for marker effects are certainly
439 more efficient for some traits and species, with the most known example being milk contents
440 in dairy cattle (VanRaden et al., 2009). This has inspired the search for its integration into
441 Single Step.

442 Bayesian regressions can be understood as inferring the variances associated to each marker
443 in the expression $Var(\mathbf{a}) = \mathbf{D}_a$, i.e. the elements $\sigma_{a,k}^2$ in the diagonal of \mathbf{D}_a being k-SNP
444 specific. Zhang et al. (2010) and Legarra et al. (2011b) checked that running a full Bayesian
445 regression to estimate breeding values, or using it to infer variances in \mathbf{D}_a to use $\mathbf{G} = \mathbf{Z}_2\mathbf{D}_a\mathbf{Z}_2'$
446 in a GBLUP gave essentially the same solution. Legarra et al. (2009) suggested to use such \mathbf{G}
447 with precomputed variances in the Single Step procedures. Makgahlela et al. (2013) picked,

448 using BayesB, either 750 or 1500 preselected markers to form $= \mathbf{Z}_2 \mathbf{D}_a \mathbf{Z}_2'$, which resulted in
449 better accuracies for milk but not for protein, and they concluded that picking the right
450 number of markers was not obvious. No other attempt has been done so far. In a similar spirit,
451 Wang et al. (2012) suggested to compute variances in \mathbf{D}_a in an iterative manner within the
452 Single Step. They obtained the marker effects from the expression
453 $\hat{\mathbf{a}}|\hat{\mathbf{u}}_2 = \mathbf{D}_a \mathbf{Z}_2' (\mathbf{Z}_2 \mathbf{D}_a \mathbf{Z}_2')^{-1} \hat{\mathbf{u}}_2$, to later infer the k-th marker variance as (proportional to) \hat{a}_k^2
454 (Sun et al., 2012). Note that this estimate is severely biased (it ignores the uncertainty in the
455 estimation of \hat{a}_k) and therefore an empirical correction needs to be applied, which is not the
456 case in true Bayesian or maximum likelihood procedures (De los Campos et al., 2009; Shen et
457 al., 2013). After computation of a new \mathbf{G} , Single Step GBLUP is rerun and markers are re-
458 estimated, and the procedure is iterated a few times. Their simulation showed an increased
459 accuracy of this method for traits with large QTLs.

460 Legarra and Ducrocq (2012) suggested two ways of dealing with Bayesian regressions. The
461 first one was to use an equivalent set of mixed model equations including marker effects:

$$\begin{pmatrix} \mathbf{X}'\mathbf{X} & \mathbf{X}'_1\mathbf{W}_1 & \mathbf{X}'_2\mathbf{W}_2\mathbf{Z}_2 \\ \mathbf{W}'_1\mathbf{X} & \mathbf{W}'_1\mathbf{W}'_1 + \mathbf{A}^{11}\lambda & \mathbf{A}^{12}\mathbf{Z}_2\lambda \\ \mathbf{Z}'_2\mathbf{W}_2\mathbf{X}_2 & \mathbf{Z}'_2\mathbf{A}^{12}\lambda & \mathbf{Z}'_2\mathbf{W}_2\mathbf{W}_2\mathbf{Z}_2 + \mathbf{Z}'_2(\mathbf{A}^{22} - \mathbf{A}^{-1}_{22})\mathbf{Z}_2\lambda + \mathbf{D}_a^{-1}\sigma_e^2 \end{pmatrix} \begin{pmatrix} \hat{\mathbf{b}} \\ \hat{\mathbf{u}}_1 \\ \hat{\mathbf{a}} \end{pmatrix} = \begin{pmatrix} \mathbf{X}'\mathbf{y} \\ \mathbf{W}'_1\mathbf{y}_1 \\ \mathbf{Z}'_2\mathbf{W}'_2\mathbf{y}_2 \end{pmatrix}$$

462

463 In this system of equations, Bayesian Regressions are accommodated by using different a
464 *priori* distributions for $\text{Var}(\mathbf{a}) = \mathbf{D}_a$ (e.g., in Bayesian Lasso the prior distribution of
465 elements in \mathbf{D}_a is double exponential). This system of equations (A1) could then be solved by
466 a Bayesian procedure such as the Gibbs sampler, which solves for \mathbf{D}_a . In the second option,
467 an equivalent iterative procedure can iterate between solutions to regular BLUP and
468 (Bayesian) genomic predictions; the results of one would be introduced into the other.
469 Because this system does not infer marker variances *per se*, it does not suffer from the bias in

470 variance estimation of Wang et al (2012). Tuning markers to be in the same scale as pedigree
 471 **in the previous set of equations or in the iterative system** would include an extra unknown for
 472 the parameter μ in Vitezica et al. (2011).

473

474 In addition, Fernando et al. (2013) recently presented another system of equations explicit on
 475 marker solutions. Equations include marker effects for *all* individuals, imputed following
 476 Gengler's method, and residual pedigree-based EBV for nongenotyped animals, ϵ . This ϵ is
 477 what remains of the breeding value after we fit (imputed) SNP effects to nongenotyped

478 individuals. Therefore total genetic value: $\mathbf{u} = \begin{pmatrix} \hat{\mathbf{Z}}_1 \\ \mathbf{Z}_2 \end{pmatrix} \mathbf{a} + \begin{pmatrix} \epsilon \\ 0 \end{pmatrix} = \hat{\mathbf{Z}} \mathbf{a} + \begin{pmatrix} \epsilon \\ 0 \end{pmatrix}$.

479 Their final Single Step mixed model equations are

$$480 \begin{pmatrix} \mathbf{X}'\mathbf{X} & \mathbf{X}'\mathbf{W}\hat{\mathbf{Z}} & \mathbf{X}'_1\mathbf{W}_1 \\ \hat{\mathbf{Z}}'\mathbf{W}'\mathbf{X} & \hat{\mathbf{Z}}'\mathbf{W}'\mathbf{W}\hat{\mathbf{Z}} + \mathbf{I} \frac{\sigma_e^2}{\sigma_a^2} & \hat{\mathbf{Z}}'_1\mathbf{W}'_1\mathbf{W}_1 \\ \mathbf{W}'_1\mathbf{X}_1 & \mathbf{W}'_1\mathbf{W}_1\hat{\mathbf{Z}}_1 & \mathbf{W}'_1\mathbf{W}_1 + \mathbf{A}^{11} \frac{\sigma_e^2}{\sigma_g^2} \end{pmatrix} \begin{pmatrix} \hat{\boldsymbol{\beta}} \\ \hat{\mathbf{a}} \\ \hat{\boldsymbol{\epsilon}} \end{pmatrix} = \begin{pmatrix} \mathbf{X}'\mathbf{y} \\ \hat{\mathbf{Z}}'\mathbf{W}'\mathbf{y} \\ \mathbf{W}'_1\mathbf{y} \end{pmatrix},$$

481 in which a Gibbs sampler can iterate to obtain Bayesian estimates. These equations are
 482 simpler than previous ones but at the cost of a very dense and large system of equations.

483 All these methods for Bayesian regressions in Single Step are largely untested, and only
 484 Wang et al. (2012) method is efficiently implemented and has been used in real data sets
 485 (Dikmen et al., 2013), for which no alternative currently exists.

486

487 4.8 Unknown parent groups

488 Missing genealogy and/or crosses are ubiquitous in animal breeding. A typical solution
 489 **consists in fitting** unknown parent groups, which model different means across groups of

490 founders well identified, i.e. belonging to different generations or breeds. BLUP equations
491 including unknown parent groups are created using an expanded inverse of the relationship
492 matrix A^{-1} (Quaas, 1988). Unfortunately, the Single Step Mixed Model equations do not
493 accommodate this well, because of the additional matrices ($G^{-1} - A_{22}^{-1}$). The problem was
494 explained in detail by Misztal et al. (2013b) who showed that proper equations would imply
495 complex terms of the form $Q_2'(G^{-1} - A_{22}^{-1})Q_2$, implying matrix Q_2 with fractions of each
496 unknown parent group for each genotyped animal. These modifications are difficult to
497 compute and program. Current alternatives involve ignoring the term (often with negligible
498 results) or using the original Westell-Robinson model, which is in the form

$$y = Xb + Qg + Wu + e$$

499 (Quaas, 1988) and fitting unknown parent groups g as covariates. This is satisfactory and
500 involves no approximations, but cumbersome to implement and of slow convergence.

501

502 *4.9 Accuracies*

503 Individual accuracies can be obtained in principle from the inverse of the Single Step mixed
504 model equations. This is impossible in practice for medium to large data sets. Therefore,
505 Misztal et al. (2013a) suggested extending known approximations in the estimation of
506 accuracy to the Single Step case. Modifications involve use of known approximations for the
507 pedigree-based BLUP and add extra information from ($G^{-1} - A_{22}^{-1}$) to each animal; then to
508 iterate the procedure. This procedure is accurate in dairy species, as attested by Misztal et al.
509 (2013a) and in Manech dairy sheep (Baloche et al., unpublished) where correlations between
510 approximate accuracies and exact accuracies from inverse of the Mixed Model Equations
511 were found to equal 0.95 in both cases.

512

513 **5. FUTURE DEVELOPMENTS**

514 Among important possible extensions, we will mention two: crosses and fit of dominance
515 effects.

516 *5.1 Crosses.*

517 Development of the Single Step has been done for purebred populations, in which heterosis is
518 absent, genetic variance is assumed constant throughout the generations, and matings are
519 (close to being) at random. In classical theory (e.g., Lo et al., 1997) populations involved in
520 crossing are assumed completely unrelated; this is subject to discussion depending on the
521 genetic architecture of the trait. For instance, Ibañez-Escriche et al. (2009) obtained the same
522 accuracy fitting markers with the same or different effects across breeds. Recently,
523 Christensen et al. (2014) presented a Single Step in these lines, where the value of a crossbred
524 animal is a sum of gametic effects, each with a different within-pure breed extended
525 relationship matrix. On the other hand, Harris and Johnson (2010, 2013) presented an
526 evaluation system for pure breeds and their complex crosses which considers different breed
527 origins but roughly the same effect of markers across breeds. These aspects need to be further
528 derived. Also, testing in real data sets is most necessary because simulations are unreliable for
529 such complex cases. However, crossbred data sets with genomic information are scarce so far.

530

531 *5.2 Dominance.*

532 Genomic predictions including dominance (e.g., Toro and Varona, 2010; Wellmann and
533 Bennewitz, 2012) are much easier than their pedigree counterparts, which are notoriously
534 difficult, in particular if inbreeding is involved (DeBoer and Hoeschele, 1993). Dominance

535 versions of GBLUP have been proposed (Su et al., 2012b; Vitezica et al., 2013) and real data
536 analysis, done (Su et al., 2012b; Ertl et al., 2013; Vitezica et al., 2013). However, these
537 methods need that genotyped animals have a phenotype, which may be precorrected. For
538 animals that have no phenotype (i.e., dairy bulls) there are no methods to generate pseudo-
539 phenotypes including dominance, because all methods to generate pseudo-data involve
540 additive relationships only. For instance, computation of DYD's in dairy cattle will average to
541 zero dominance deviations of the offspring. Therefore Single Step methods for dominance are
542 highly relevant, yet a simple combination of pedigree-based and marker-based methods is
543 difficult because the pedigree-based method is already difficult.

544

545 **6. OBSCURE POINTS AND LIMITS**

546 *6.1 Treatment of linkage.*

547 Markers are physically linked and their co-occurrence is correlated. However, most genomic
548 prediction models, including Bayesian Regressions and the Single Step, assume markers to be
549 unlinked. In addition, the pedigree-based matrix \mathbf{A} assumes loci as unlinked as well.
550 Meuwissen et al. (2011) suggested a modified \mathbf{H} matrix in which pedigree relationships
551 would not be included using pedigree relationships \mathbf{A} , but using \mathbf{G}_{FG} , the Fernando and
552 Grossmann (1989) covariance matrix using pedigree and markers. The latter would be
553 computed by means of iterative peeling, producing relationships for all individuals, genotyped
554 or not. This procedure provides in principle a more accurate relationship matrix, and therefore
555 should result in more accurate Single Step evaluations. However, the extent of this extra
556 accuracy has not been evaluated in realistic simulations (e.g., with large genomes and large
557 number of animals) or in real life data sets and it is unknown how this method scales to large
558 pedigrees.

559 *6.2 Convergence of solvers.*

560 The convergence rate with regular Single Step when solved by PCG iteration depends on
561 species. The rate is similar to BLUP and poses no problem with complete pedigree and a
562 uniform base population (e.g., chicken). The rate is also good with high-accuracy genotyped
563 animals (dairy bulls). The rate can be poor with complex models when the pedigree contains
564 many generations of animals without phenotypes. In such a case, restricting the pedigree to
565 fewer old animals improves the rate. Poor convergence rate in some models is due to
566 incompatibility between \mathbf{G} and \mathbf{A}_{22} when the pedigree has missing animals across generations
567 (Misztal et al., 2013). When \mathbf{G} is scaled for an average \mathbf{A}_{22} , elements of \mathbf{A}_{22}^{-1} due to animals
568 with very long pedigree are larger. Solutions to this problem include modifications to \mathbf{A} (e.g.,
569 as in Christensen, 2012), or pedigree or even phenotype truncations. Lourenco et al. (2014)
570 investigated the effect of cutting pedigrees and phenotypes on accuracy for the youngest
571 generation. Use of data beyond 2 generations of phenotypes and 4 generations of pedigree did
572 not improve the accuracy while increasing computing costs.

573 In large data sets with many genotyped individuals (e.g., with genotyped cows) there are
574 reports of lack of, or very slow, convergence (Harris et al., 2013; VanRaden, unpublished).
575 This raises the question if the typical form of the mixed model equations for single-Step,
576 including \mathbf{G} and \mathbf{A}_{22} is the most appropriate, or alternative forms based on marker effects such
577 as those presented by Legarra and Ducrocq (2012) or Fernando et al. (2013) are better
578 numerically conditioned. No real data testing of these approaches has been shown so far. A
579 limit to testing these approaches is the availability of very general software for BLUP.
580 General software (multiple trait, multiple effects, etc.) does not exist for marker-based
581 methods.

582 *6.3 Computational limits.*

583 Computing and inverting \mathbf{G} and \mathbf{A}_{22} is challenging and of cubic cost, which will
584 eventually preclude its use for, say, >100,000 animals, and alternatives have been suggested
585 (Legarra and Ducrocq, 2012; Fernando et al., 2013) but not thoroughly tested. These
586 alternatives would be either highly parallelizable or use indirect representations avoiding
587 explicit computations. However, so far, problems of convergence seem more limiting than
588 size.

589

590 **7. CURRENT STATE AND PRACTICAL EXPERIENCES**

591 *7.1 Dairy sheep.*

592 In France, the Lacaune, Manech and Basco-Bernaise genomic evaluations use Single Step in
593 its typical form, with corrections of \mathbf{G} to match \mathbf{A}_{22} and with the fit of unknown parent groups
594 as covariates. Preliminary research did not show an added accuracy of **Bayesian Regressions**
595 (Duchemin et al., 2012). Single step results in higher accuracy than GBLUP with pseudo-
596 phenotypes (Baloche et al., 2014) and in a much simpler implementation. Single Step will be
597 the method for genomic prediction in the future Lacaune dairy sheep genomic selection
598 scheme.

599 *7.2 Dairy goat.*

600 In France, the dairy goat population is testing genomic selection procedures with the Single
601 Step as the evaluation tool (Carillier et al., 2013) although it is very soon to establish its
602 impact.

603 *7.3 Pigs.*

604 In Denmark, routine genetic evaluation of the three DanBred breeds Duroc, Landrace and
605 Yorkshire has since October 2011 been made by Single-Step in its typical form, with
606 corrections of \mathbf{G} to match \mathbf{A}_{22} . The implementation of genomic evaluation via Single-Step
607 was straight-forward and it has resulted in increased accuracy compared to the traditional
608 genetic evaluation. Breeding companies PIC and ToPigs also use Single Step for genomic
609 predictions.

610 *7.4 Dairy cattle.*

611 National evaluations are based on multiple step procedures, but most countries are willing to
612 change to Single Step, and many are experimenting (e.g., VanRaden, unpublished; Koivula et
613 al., 2012; Harris et al., 2013). The reason for this change is the conceptual and practical
614 simplicity of the Single Step, and its ability to account for genomic preselection (Patry and
615 Ducrocq, 2011). Due to abundance of data and completeness of genotyping, tests show
616 equivalent accuracies of Single Step and multiple step procedures (e.g., Aguilar et al., 2010).
617 SsGBLUP was always more accurate than GBLUP for several milkability traits (Gray et al.,
618 2012), and slightly more accurate for test-day models (Koivula et al., 2012). Also, Pribyl et al.
619 (2013) showed higher accuracy of the Single Step for Check Republic data.

620 *7.5 Beef cattle*

621 There are no studies on the application of Single Step to real data sets. These data sets are
622 more complex for genomic evaluation than other species because of missing relationships,
623 smaller sibships, and the presence of maternal effects. Real data studies are therefore much
624 needed. However, in a simulation study by Lourenco et al. (2013), accuracies of genomic
625 predictions with ssGBLUP were always higher than with BLUP, which was not the case with
626 BayesC. This was particularly true for maternal traits.

627 *7.6 Chicken*

628 In studies on decay of genomic prediction over generations (Wolc et al., 2011), BayesB was
629 more accurate than single-trait GBLUP but less accurate than 2-trait GBLUP; in that study,
630 GBLUP was applied to a reduced animal model and was equivalent to ssGBLUP. Chen et al.
631 (2011a,b) also showed higher accuracies of Single Step than with Bayesian regressions.

632

633 **8. SOFTWARE**

634 To our knowledge, the only publicly available software packages which can directly run
635 Single Step evaluations are the BLUPF90 family of programs (Misztal et al., 2002;
636 <http://nce.ads.uga.edu/wiki>) and software DMU (Madsen and Jensen, 2012,
637 <http://www.dmu.agrsci.dk/>) in which it is fully implemented including regular BLUP, REML,
638 Gibbs samplers (only BLUPF90), threshold models, generalized linear mixed models (only
639 DMU) and iteration on data for very large data sets, and several options (most of them
640 mentioned above). Software Mix99 (Vuori et al., 2006) has been modified to include Single
641 Step, although these modifications are not publicly available. Public packages such as
642 Wombat (Meyer, 2013; <http://didgeridoo.une.edu.au/km/wombat.php>) or ASREML
643 (<http://www.vsnr.co.uk/software/asreml>) can include covariance matrices computed
644 externally, and therefore matrix H^{-1} needs to be computed with an external tool and then fit
645 into the model.

646

647 **9. CONCLUSION: OVERALL BENEFITS AND DRAWBACKS OF THE SINGLE** 648 **STEP**

649 The Single Step provides a simple method to combine all information in a simple manner,
650 with the additional advantage of requiring little changes to existing software. Accuracy is

651 usually as high as, if not greater than, any other method. Some studies concerning accuracy of
652 the Single Step have been gathered in Table 1. Beyond its extra accuracy, it has the following
653 interesting properties:

- 654 1. Automatic accounting of all relatives of genotyped individuals and their performances.
- 655 2. Simultaneous fit of genomic information and estimates of other effects (e.g.,
656 contemporary groups). Therefore not loss of information.
- 657 3. Feedback: the extra accuracy in genotyped individuals is transmitted to all their
658 relatives (*e.g.* Christensen et al., 2012).
- 659 4. Simple extensions. Because this is a linear BLUP-like estimator, the extension to more
660 complicated models (multiple trait, threshold traits, test day records) is immediate.
661 Any model fit using relationship matrices can be fit using combined relationship
662 matrices.
- 663 5. Analytical framework. The Single Step provides an analytical framework for further
664 developments. This is notoriously difficult with pseudo-data.

665 As drawbacks, one can cite the following:

- 666 1. Programming complexity to fit complicated models for marker effects (Bayesian
667 Regressions, machine learning algorithms, etc.).
- 668 2. Lack of experience on very large data sets.
- 669 3. Long computing times with current Single Step algorithms methods, for very large
670 data sets.
- 671 4. Lack of an easy and elegant way of considering major genes in a multiple trait setting,
672 this is a drawback of multiple step methods as well.

673

674 TABLE 1 HERE

675

676

677 REFERENCES

- 678 Abraham, K.J., Totir, L.R., Fernando, R.L., 2007. Improved techniques for sampling complex
679 pedigrees with the Gibbs sampler. *Gen Sel Evol* 39, 27-38.
- 680 Aguilar, I., Misztal, I., Johnson, D.L., Legarra, A., Tsuruta, S., Lawlor, T.J., 2010. Hot topic:
681 a unified approach to utilize phenotypic, full pedigree, and genomic information for genetic
682 evaluation of Holstein final score. *J Dairy Sci* 93, 743-752.
- 683 Aguilar, I., Misztal, I., Legarra, A., Tsuruta, S., 2011. Efficient computations of genomic
684 relationship matrix and other matrices used in the single-step evaluation. *J Anim Breed Genet*
685 128, 422-428.
- 686 Aguilar, I., Legarra, A., Tsuruta, S., Misztal, I., 2013. Genetic evaluation using unsymmetric
687 single step genomic methodology with large number of genotypes. *Interbull Bulletin* 47.
- 688 Baloche, G., Legarra, A., Sallé, G., Larroque, H., Astruc, J.M., Robert-Granié, C., Barillet, F.,
689 2014. Assessment of accuracy of genomic prediction for French Lacaune dairy sheep. *J Dairy*
690 *Sci* 97, 1107-1116.
- 691 Carillier, C., Larroque, H., Palhière, I., Clément, V., Rupp, R., Robert-Granié, C., 2013. A
692 first step toward genomic selection in the multi-breed French dairy goat population. *J Dairy*
693 *Sci* 96, 7294-7305.
- 694 Casella, G., Berger, R.L., 1990. *Statistical inference*. Duxbury Press Belmont, CA.
- 695 Chen, C., Misztal, I., Aguilar, I., Tsuruta, S., Aggrey, S., Wing, T., Muir, W., 2011a.
696 Genome-wide marker-assisted selection combining all pedigree phenotypic information with
697 genotypic data in one step: An example using broiler chickens. *J Anim Sci* 89, 23-28.
- 698 Chen, C.Y., Misztal, I., Aguilar, I., Legarra, A., Muir, W.M., 2011b. Effect of different
699 genomic relationship matrices on accuracy and scale. *J Anim Sci* 89, 2673-2679.
- 700 Christensen, O.F., 2012. Compatibility of pedigree-based and marker-based relationship
701 matrices for single-step genetic evaluation. *Gen Sel Evol* 44, 37.

702 Christensen, O.F., Lund, M.S., 2010. Genomic prediction when some animals are not
703 genotyped. *Gen Sel Evol* 42, 2.

704 Christensen, O., Madsen, P., Nielsen, B., Ostersen, T., Su, G., 2012. Single-step methods for
705 genomic evaluation in pigs. *Animal* 6, 1565-1571.

706 Christensen, O.F., Madsen, P., Nielsen, B. Su, G., 2014. Genomic evaluation of both purebred
707 and crossbred performances. *Gen Sel Evol* 46, 23

708 Cochran, W., 1951. Improvement by means of selection. Second Berkeley Symposium on
709 Mathematical Statistics and Probability, pp. 449-470.

710 Cockerham, C.C., 1969. Variance of gene frequencies. *Evolution* 23, 72-84.

711 Cockerham, C.C., 1973. Analyses of gene frequencies. *Genetics* 74, 679.

712 Colleau, J.J., 2002. An indirect approach to the extensive calculation of relationship
713 coefficients. *Gen Sel Evol* 34, 409-422.

714 De Boer, I., Hoeschele, I., 1993. Genetic evaluation methods for populations with dominance
715 and inbreeding. *Theor Appl Gen* 86, 245-258.

716 de los Campos, G., Naya, H., Gianola, D., Crossa, J., Legarra, A., Manfredi, E., Weigel, K.,
717 Cotes, J.M., 2009. Predicting quantitative traits with regression models for dense molecular
718 markers and pedigree. *Genetics* 182, 375-385.

719 Dikmen, S., Cole, J.B., Null, D.J., Hansen, P.J., 2013. Genome-wide association mapping for
720 identification of quantitative trait loci for rectal temperature during heat stress in Holstein
721 cattle. *PLoS ONE* 8, e69202.

722 Duchemin, S., Colombani, C., Legarra, A., Baloche, G., Larroque, H., Astruc, J.-M., Barillet,
723 F., Robert-Granié, C., Manfredi, E., 2012. Genomic selection in the French Lacaune dairy
724 sheep breed. *J Dairy Sci* 95, 2723-2733.

725 Ertl, J., Legarra, A., Vitezica, Z.G., Varona, L., Edel, C., Reiner, E., Götz, K.-U., 2013.
726 Genomic analysis of dominance effects in milk production and conformation traits of

727 Fleckvieh cattle. *Interbull Bulletin* 47.

728 Fernando, R., Gianola, D., 1986. Optimal properties of the conditional mean as a selection
729 criterion. *Theor Appl Gen* 72, 822-825.

730 Fernando, R.L., Grossman, M., 1989. Marker assisted prediction using best linear unbiased
731 prediction. *Gen Sel Evol* 21, 467-477.

732 Fernando, R.L., Garrick, D.J., Dekkers, J.C.M., 2013. Bayesian regression method for
733 genomic analyses with incomplete genotype data. European Federation of Animal Science.
734 Wageningen Press, Nantes, France, p. 225.

735 Forni, S., Aguilar, I., Misztal, I., 2011. Different genomic relationship matrices for single-step
736 analysis using phenotypic, pedigree and genomic information. *Gen Sel Evol* 43, 1.

737 Garcia-Cortes, L.A., Legarra, A., Chevalet, C., Toro, M.A., 2013. Variance and Covariance of
738 Actual Relationships between Relatives at One Locus. *PLoS ONE* 8, e57003.

739 Garrick, D.J., Taylor, J.F., Fernando, R.L., 2009. Deregressing estimated breeding values and
740 weighting information for genomic regression analyses. *Gen Sel Evol* 41, 44.

741 Gengler, N., Mayeres, P., Szydlowski, M., 2007. A simple method to approximate gene
742 content in large pedigree populations: application to the myostatin gene in dual-purpose
743 Belgian Blue cattle. *animal* 1, 21-28.

744 Gianola, D., Fernando, R.L., 1986. Bayesian Methods in Animal Breeding Theory. *J Anim*
745 *Sci* 63, 217.

746 Gray, K.A., Cassady, J.P., Huang, Y., Maltecca, C., 2012. Effectiveness of genomic
747 prediction on milk flow traits in dairy cattle. *Gen Sel Evol* 44:24

748 Harris, B.L., Johnson, D.L., 2010. Genomic predictions for New Zealand dairy bulls and
749 integration with national genetic evaluation. *J Dairy Sci* 93:1243-1252.

750 Harris, B.L., Winkelman, A.M., Johnson, D.L., 2013. Impact of including a large number of
751 female genotypes on genomic selection. *Interbull Bulletin* 47.

752 Hayes, B.J., Visscher, P.M., Goddard, M.E., 2009. Increased accuracy of artificial selection
753 by using the realized relationship matrix. *Genet Res* 91, 47-60.

754 Henderson, C.R., 1973. Sire evaluation and genetic trends. In *Proceedings of the animal*
755 *breeding and genetics symposium in honor of Dr. Jay L. Lush* pp. 10-41.

756 Hickey, J.M., Kinghorn, B.P., Tier, B., van der Werf, J.H., Cleveland, M.A., 2012. A phasing
757 and imputation method for pedigreed populations that results in a single-stage genomic
758 evaluation. *Gen Sel Evol* 44.

759 Hill, W.G., Weir, B.S., 2011. Variation in actual relationship as a consequence of Mendelian
760 sampling and linkage. *Genet Res (Camb)*, 1-18.

761 Ibáñez-Escriche, N., Fernando, R.L., Toosi, A., Dekkers, J.C.M., 2009. Genomic selection of
762 purebreds for crossbred performance. *Gen Sel Evol* 41, 12.

763 Koivula, M., Strandén, I., Pösö, J., Aamand, G.P., Mäntysaari, E.A., 2012. Single step
764 genomic evaluations for the Nordic Red Dairy cattle test day data. *Interbull Bull*, 46.

765 Legarra, A., Misztal, I., 2008. Technical note: Computing strategies in genome-wide
766 selection. *J Dairy Sci* 91, 360-366.

767 Legarra, A., Robert-Granié, C., Manfredi, E., Elsen, J.-M., 2008. Performance of genomic
768 selection in mice. *Genetics* 180, 611-618.

769 Legarra, A., Aguilar, I., Misztal, I., 2009. A relationship matrix including full pedigree and
770 genomic information. *J Dairy Sci* 92, 4656-4663.

771 Legarra, A., Calenge, F., Mariani, P., Velge, P., Beaumont, C., 2011a. Use of a reduced set of
772 single nucleotide polymorphisms for genetic evaluation of resistance to *Salmonella* carrier
773 state in laying hens. *Poultry Sci*, 90, 731-736.

774 Legarra, A., Robert-Granié, C., Croiseau, P., Guillaume, F., Fritz, S., 2011b. Improved Lasso
775 for genomic selection. *Genet Res (Camb)* 93, 77-87.

776 Legarra, A., Ducrocq, V., 2012. Computational strategies for national integration of

777 phenotypic, genomic, and pedigree data in a single-step best linear unbiased prediction. J
778 Dairy Sci 95, 4629-4645.

779 Lo, L., Fernando, R., Grossman, M., 1997. Genetic evaluation by BLUP in two-breed
780 terminal crossbreeding systems under dominance. J Anim Sci 75, 2877-2884.

781 Lourenco, D., Misztal, I., Wang, H., Aguilar, I., Tsuruta, S., Bertrand, J., 2013. Prediction
782 accuracy for a simulated maternally affected trait of beef cattle using different genomic
783 evaluation models. J Anim Sci 91, 4090-4098.

784 Lourenco, D., Misztal, I., Tsuruta, S., Aguilar, I., Lawlor, T. J., Forni, S., Weller, J. I. 2014.
785 Are evaluations on young genotyped animals benefiting from the past generations? J Dairy
786 Sci, in press.

787 Madsen, P., Jensen, J., 2000. A user's guide to DMU. A package for analysing multivariate
788 mixed models. Version 6, 1-33.

789 Makgahlela, M. L., Knürr, T., Aamand, G., Strandén, I., Mäntyselä, E., 2013. Single step
790 evaluations using haplotype segments. Interbull Bulletin, 47.

791 Meuwissen, T.H.E., Hayes, B.J., Goddard, M.E., 2001. Prediction of total genetic value using
792 genome-wide dense marker maps. Genetics 157, 1819-1829.

793 Meuwissen, T., Luan, T., Woolliams, J., 2011. The unified approach to the use of genomic
794 and pedigree information in genomic evaluations revisited. J Anim Breed Genet 128, 429-
795 439.

796 Meyer, K., 2007. WOMBAT—A tool for mixed model analyses in quantitative genetics by
797 restricted maximum likelihood (REML). Journal of Zhejiang University Science B, 8(11),
798 815-821.

799 Misztal, I., Tsuruta, S., Strabel, T., Auvray, B., Druet, T., and Lee, D. H. (2002). BLUPF90
800 and related programs (BGF90). In Proceedings of the 7th World Congress on Genetics
801 Applied to Livestock Production, Montpellier, France, August, 2002. Session 28. (pp. 1-2).

802 Institut National de la Recherche Agronomique (INRA).

803 Misztal, I., Legarra, A., Aguilar, I., 2009. Computing procedures for genetic evaluation
804 including phenotypic, full pedigree, and genomic information. *J Dairy Sci* 92, 4648-4655.

805 Misztal, I., Tsuruta, S., Aguilar, I., Legarra, A., VanRaden, P., Lawlor, T., 2013a. Methods to
806 approximate reliabilities in single-step genomic evaluation. *J Dairy Sci.* 96, 647-654

807 Misztal, I., Vitezica, Z., Legarra, A., Aguilar, I., Swan, A., 2013b. Unknown-parent groups in
808 single-step genomic evaluation. *J Anim Breed Genet.* 130, 252-258.

809 Patry, C., Ducrocq, V., 2011. Evidence of biases in genetic evaluations due to genomic
810 preselection in dairy cattle. *J Dairy Sci* 94, 1011-1020.

811 Powell, J.E., Visscher, P.M., Goddard, M.E., 2010. Reconciling the analysis of IBD and IBS
812 in complex trait studies. *Nat Rev Genet* 11, 800-805.

813 Přibyl, J., Madsen, P., Bauer, J., Přibyl, J., Šimečková, M., Vostrý, L., Zavadilová, L.,
814 2013. Contribution of domestic production records, Interbull estimated breeding values, and
815 single nucleotide polymorphism genetic markers to the single-step genomic evaluation of
816 milk production. *J Dairy Sci* 96, 1865-1873.

817 Quaas, R.L., 1988. Additive genetic model with groups and relationships. *J Dairy Sci* 71,
818 1338-1345.

819 Ricard, A., Danvy, S., Legarra, A., 2013. Computation of deregressed proofs for genomic
820 selection when own phenotypes exist with an application in French show-jumping horses. *J*
821 *Anim Sci* 91, 1076-1085.

822 Ritland, K., 1996. Estimators for pairwise relatedness and individual inbreeding coefficients.
823 *Genetical Research* 67, 175-185.

824 Shen, X., Alam, M., Fikse, F., Rönnegård, L., 2013. A novel generalized ridge regression
825 method for quantitative genetics. *Genetics* 193, 1255-1268.

826 Smith, H.F., 1936. A discriminant function for plant selection. *Annals of Eugenics* 7, 240-

827 250.

828 Strandén, I., Lidauer, M., 1999. Solving large mixed linear models using preconditioned
829 conjugate gradient iteration. *J Dairy Sci* 82, 2779-2787.

830 Strandén, I., Garrick, D.J., 2009. Technical note: Derivation of equivalent computing
831 algorithms for genomic predictions and reliabilities of animal merit. *J Dairy Sci* 92, 2971-
832 2975.

833 Su, G., Christensen, O.F., Ostersen, T., Henryon, M., Lund, M.S., 2012a. Estimating additive
834 and non-additive genetic variances and predicting genetic merits using genome-wide dense
835 single nucleotide polymorphism markers. *PLoS ONE* 7, e45293.

836 Su, G., Madsen, P., Nielsen, U.S., Mäntysaari, E.A., Aamand, G.P., Christensen, O.F., Lund,
837 M.S., 2012b. Genomic prediction for Nordic Red Cattle using one-step and selection index
838 blending. *J Dairy Sci* 95, 909-917.

839 Sun, X., Qu, L., Garrick, D.J., Dekkers, J.C., Fernando, R.L., 2012. A Fast EM Algorithm for
840 BayesA-Like Prediction of Genomic Breeding Values. *PLoS ONE* 7, e49157.

841 Sun, C., Van Raden, P., 2013. Mating programs including genomic relationships. *J Dairy Sci*
842 96, 653.

843 Tanner, M.A., Wong, W.H., 1987. The calculation of posterior distributions by data
844 augmentation. *J Amer Stat Assoc* 82, 528-540.

845 Toro, M.Á., García-Cortés, L.A., Legarra, A., 2011. A note on the rationale for estimating
846 genealogical coancestry from molecular markers. *Gen Sel Evol* 43, 27.

847 Toro, M.A., Varona, L., 2010. A note on mate allocation for dominance handling in genomic
848 selection. *Gen Sel Evol* 42, 33.

849 VanRaden, P.M., 2008. Efficient Methods to Compute Genomic Predictions. *J. Dairy Sci.* 91,
850 4414-4423.

851 VanRaden, P., Wiggans, G., 1991. Derivation, calculation, and use of national animal model

852 information. *J Dairy Sci* 74, 2737-2746.

853 VanRaden, P.M., Tassell, C.P.V., Wiggans, G.R., Sonstegard, T.S., Schnabel, R.D., Taylor,
854 J.F., Schenkel, F.S., 2009. Invited review: reliability of genomic predictions for North
855 American Holstein bulls. *J Dairy Sci* 92, 16-24.

856 VanRaden, P., Null, D., Sargolzaei, M., Wiggans, G., Tooker, M., Cole, J., Sonstegard, T.,
857 Connor, E., Winters, M., van Kaam, J., 2013. Genomic imputation and evaluation using high-
858 density Holstein genotypes. *J Dairy Sci* 96, 668-678.

859 Vitezica, Z., Aguilar, I., Misztal, I., Legarra, A., 2011. Bias in genomic predictions for
860 populations under selection. *Genetics Research* 93, 357-366.

861 Vitezica, Z.G., Varona, L., Legarra, A., 2013. On the additive and dominant variance and
862 covariance of individuals within the genomic selection scope. *Genetics* 195, 1223-1230.

863 Vuori, K., Strandén, I., Lidauer, M., Mäntysaari, E., 2006. MiX99-effective solver for large
864 and complex linear mixed models. *Proceedings of the 8th World Congress on Genetics
865 Applied to Livestock Production, Belo Horizonte, Minas Gerais, Brazil, 13-18 August, 2006.
866 Instituto Prociência, pp. 27-33.*

867 Wang, H., Misztal, I., Aguilar, I., Legarra, A., Muir, W., 2012. Genome-wide association
868 mapping including phenotypes from relatives without genotypes. *Genetics Research* 94, 73-
869 83.

870 Wellmann, R., Bennewitz, J., 2012. Bayesian models with dominance effects for genomic
871 evaluation of quantitative traits. *Genetics Research* 94, 21.

872 Wolc, A., Stricker, C., Arango, J., Settar, P., Fulton, J.E., O'Sullivan, N.P., Preisinger, R.,
873 Habier, D., Fernando, R., Garrick, D.J., Lamont, S.J., Dekkers, J.C.M., 2011. Breeding value
874 prediction for production traits in layer chickens using pedigree or genomic relationships in a
875 reduced animal model. *Gen Sel Evol* 43, 5.

876 Zhang, Z., Liu, J., Ding, X., Bijma, P., de Koning, D.-J., Zhang, Q., 2010. Best linear

877 unbiased prediction of genomic breeding values using a trait-specific marker-derived
878 relationship matrix. PLoS ONE 5, e12648.

879

880

881 **Table 1. Accuracy of Single Step versus other methods in some species**

Authors	Single Step	Multiple step	Pedigree BLUP	Species, trait
Aguilar et al., 2010	0.70	0.70	0.60	Dairy cattle, final score
Baloche et al., 2013	0.47	0.43	0.32	Milk yield, dairy sheep
Chen et al., 2011b*	0.36		0.20	Breast meat, chicken
Chen et al., 2011a	0.37	0.09	0.28	Leg Score, chicken
Christensen et al., 2012*	0.35	0.35	0.18	Daily gain, pigs
Aguilar et al., 2011	0.39		0.26	Conception rate at first parity

882 *predictive abilities: $r(y, \hat{u})$

883

884

885

886