

## 1 **Bayesian genomic prediction of junctional epidermolysis bullosa in sheep**

2 **Abstract:** Junctional epidermolysis bullosa (JEP) is a heritable skin and mucosa disorders  
3 condition in association with mendelian mutations in sheep. The purpose of this investigation  
4 is to explore the relationship between different priors, linkage disequilibrium and single  
5 nucleotide polymorphisms (SNPs) selection methods to accuracy of Bayesian GP of JEP in  
6 sheep. 92 Spanish Churra sheep breed genotyped by 40668 SNP markers. Bayes  $C\pi$  shown to  
7 have slightly higher predicted accuracy [0.724 (0.113)] by unselected data. Prediction  
8 performance of the Bayesian GP models was found to be similar after correction for LD. There  
9 was a significant difference between predicted accuracies due to the SNPs selection by ranked  
10 p values of whole and training only dataset using linear model. The relevance of genetic  
11 architecture in conjugate to the prior distributions clearly supported by the unselected data. The  
12 most obvious finding emerge from this study is that preselection of SNPs referring to genetic  
13 architecture of the phenotype may lower the needs of computational load.

14 **Key words:** Bayesian models, genomic prediction, junctional epidermolysis bullosa

### 15 **1. Introduction**

16 Genomics have emerged as powerful platform in animal breeding and genetics due to decreased  
17 costs of molecular markers [1]. Genomic prediction (GP) is important for a wide range of  
18 scientific and industrial process in animal breeding including: detection of genes in connection  
19 with phenotypes and prediction of genomic breeding values [2, 3]. The main challenge faced  
20 by many researchers is the GP of disease statutes of the animals using single nucleotide  
21 polymorphisms (SNPs) to obtain clinical diagnostic systems [4].

22 Scholars have debated the impact of genetic architecture of the phenotype [5], preselection of  
23 SNPs [4], and differences between GP methods [4] for explaining the variation in GP accuracy.

24 In the literature of Bayesian GP, the relative importance of prior distributions has been subject  
25 to considerable discussion [6]. Investigating genetic architecture of the phenotypes in terms of

1 prior distributions is a continuing concern to obtain higher GP accuracies. It has conclusively  
2 been shown that different priors lead to different genetic architectures in terms of number and  
3 action of the genes and level of linkage disequilibrium (LD) [7]. The GP research to date has  
4 tended to focus on quantitative phenotypes rather than binary disease statuses.

5 Junctional epidermolysis bullosa (JEP) is a heritable skin and mucosa disorders condition in  
6 association with mendelian mutations in sheep [8, 9]. However genetic factors found to be  
7 influencing epidermolysis bullosa have been explored in several organisms including cattle  
8 [10], sheep [8, 11, 12], horse [13], dog, cats and rats [13]. Surveys in mammals such as that  
9 conducted by [10] have shown that hundreds of mutations in association with 18 genes have  
10 been molecularly characterized for epidermolysis bullosa. The purpose of this investigation is  
11 to explore the relationship between different priors, LD and SNPs selection methods to accuracy  
12 of Bayesian GP of JEP in sheep.

## 13 **2. Materials and methods**

14 92 (17 cases and 75 controls) Spanish Churra sheep breed genotyped by 40668 SNP markers.  
15 Phenotypes were assessed by visual inspection of the sheeps and recorded as a binary trait.  
16 More details about the dataset could be found at [8]. SNPs were analysed by PLINK [14] for  
17 quality control based on minor allele frequencies ( $<0.05$ ), calling rate of SNPs ( $>0.90$ ), Hardy-  
18 Weinberg proportions ( $P<1E-07$ ), and optionally Linkage Disequilibrium (LD) ( $r^2>0.7$ ).

19 The evolution of the Bayesian GP models was based on cross validation of the genotypic and  
20 phenotypic datasets over training and testing partions. Splitting the data as training (%80 of  
21 animals) and testing (% 20 of animals) are common for evolution of GP methods [4, 5]. Area  
22 under the curve (AUC) approach was used to obtain the accuracies over training and testing  
23 partions by using 10 fold cross validations of Bayesian GP methods as was defined in [4].

1 Bayesian ridge regression (BRR), Bayesian (Least Absolute Shrinkage and Selection Operator)  
2 LASSO (BL), Bayes A, Bayes B and Bayes  $C\pi$  [7] are currently the most popular Bayesian GP  
3 methods for investigating animal breeding datasets. To obtain GP for animals  $y = \sum_i^n X_i \hat{g}_i$   
4 could be used where  $y$  is the phenotype (1 for case of JEP, 0 for healthy control),  $n$  is the number  
5 of SNPs,  $X_i$  is a design matrix connecting animals to genotypes at SNP  $i$ , and  $\hat{g}_i$  is the predicted  
6 effect of the genotype at SNP  $i$ .  $\hat{g}_i$  have been used to investigate the genetical properties on the  
7 phenotype with referring various prior assumptions regarding genetic architecture of the  
8 phenotype. BRR [15] assumes same additive genetic variance for all SNPs by using normal  
9 prior distribution. BL [16] assumes Laplace prior distribution for shrinking many of the SNPs  
10 towards to zero. Bayes A [15] assumes for the distribution of SNP effects is the Student's  $t$   
11 distribution. However, there are certain drawbacks associated with the use of Bayes A including  
12 non-zero SNP effects over genome. The use of mixture models has a relatively long tradition  
13 within GP [17]. One advantage of Bayes B [15] is that it avoids the problem of non-zero SNPs  
14 effects by using mixture of two prior distributions:

15  $\sigma_{gi}^2 = 0$  with probability  $\pi$ ,

16  $\sigma_{gi}^2 \sim \chi^{-2}(v, S)$  with probability  $(1-\pi)$

17 where  $\sigma_{gi}^2$  is the additive genetic variance of SNP  $i$ , with  $v=4.234$ ,  $S=0.0429$  [15],  $\pi$  (assumed  
18 to be 0.5) is the probability that the SNP has no effect on the phenotype. One possible  
19 improvement over Bayes B could be obtained by predicting  $\pi$  parameter in Bayes  $C\pi$ . Bayesian  
20 GP analysed by the BGLR package [18] with 52000 markov chain monte carlo iterations by  
21 6000 burn-in period.

22 The literature on preselection of SNPs for GP has revealed the emergence of several contrasting  
23 themes [4, 19] referring genetic architecture of the phenotypes. SNPs were filtered for the  
24 stratified training and testing GWAS results set by P values ( $<0.05$ ) and full data set ranked P

1 values ( $<0.05$ ) of linear mixed model (LMM) [20]. Different from Bayes  $C\pi$  : BayesR assumes  
2 prior distributions with four mixture components of Gaussian distribution to model SNPs  
3 effects. LMM used a single SNP regression model hence only SNPs with large effects could be  
4 detected.

### 5 **3. Results and Discussion**

6 After quality control process: 40642 SNPs with 92 sheep obtained. After filtering out highly  
7 correlated SNPs from the genotypic dataset, 25254 SNPs obtained (Figure 1). Whole and  
8 training only LMM detected 2401 SNPs and 2120 SNPs in association with JEB respectively.  
9 Figure 2 and Table 1 compares the prediction accuracies obtained from Bayesian GP models  
10 under different experimental designs. This Table 1 is quite revealing in several ways. First,  
11 Bayes  $C\pi$  shown to have slightly higher predicted accuracy by unselected data. Prediction  
12 performance of the Bayesian GP models was found to be similar after correction for LD. There  
13 was a significant difference between predicted accuracies due to the SNPs selection by whole  
14 and training only LMM analyses. Interestingly, the full LMM based preselected data gave the  
15 highest GP accuracy with relatively smaller sampling size and smaller standard errors (Figure  
16 2) compared with the other experimental designs in Table 1. Smallest sampling size (2120  
17 SNPs) with relatively higher prediction accuracies were obtained by the preselection of SNPs  
18 using training only LMM. Prediction accuracies was found to be similar over different Bayesian  
19 GP models in each SNPs selection methods.

20 This study set out with the aim of assessing the importance of Bayesian GP of JEP in sheep  
21 under different experimental settings. The results of this study indicate that it is possible to  
22 predict JEP in sheep using SNPs data. Increased GP accuracy over LMM selected data (Table  
23 1) in this study corroborates with hypothesis of mendelian inheritance pattern for JEP [9]. LMM  
24 provided the SNPs with largest set of mendelian effects due to its single SNP regression

1 algorithm. This finding broadly supports the work of other studies in this area linking mendelian  
2 mutations with JEP [8, 13]. Increased GP accuracy over, Bayes  $C\pi$  (in the unselected dataset)  
3 in this study corroborates these earlier findings of mendelian inheritance of JEP [8].

4 [21] demonstrated that postulated prior values for Bayesian GP became more important  
5 especially with small datasets. Surprisingly, no differences were found in Bayesian prediction  
6 accuracies over preselected datasets using different prior distributions (Table 1). These results  
7 could be explained by the fact that, filtering the SNPs for high LD or preselection of the SNPs  
8 for their effect sizes (by LMM) reduced the genotypic variation in the data and all Bayesian  
9 models started to give similar accuracies. These results reflect those of [22] who also found the  
10 flexibility and interpretability of GP obtained by BayesR under various simulation experiments  
11 including corrections for LD blocks. In accordance with the present results [4] reported their  
12 Bayesian GP accuracies were found to be similar due to unmatched genetic architecture of the  
13 phenotype and postulated prior distributions.

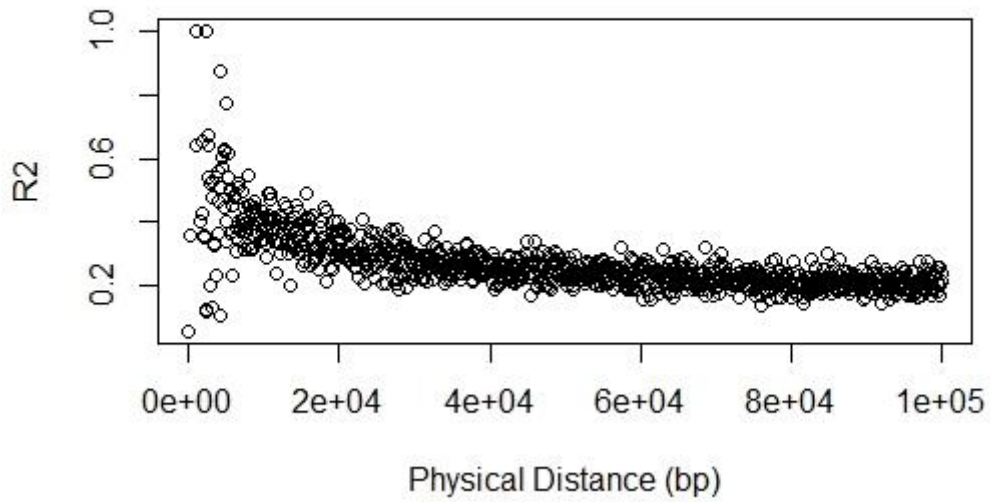
14 Including preselected SNPs in the model improved prediction accuracies using GWAS results  
15 of whole genotypic dataset (Table 1). Preselection of SNPs beneficial not only for improvement  
16 of prediction accuracies but also for reduction of dimension of the genotypic dataset [23].  
17 Instead of 40642 SNPs: LMM SNPs selection models GPs obtained by 2401 SNPs. This finding  
18 is consistent with that of [24] who reported advantages of using BayesR in GP for take into  
19 account of variants with large effects. It is possible to hypothesises that employing SNPs selection  
20 models in conjugate with genetic architecture of the phenotype would be more efficient  
21 compared with other data reduction techniques as such principal component analyses. In  
22 accordance with the present results [25] point out that Bayesian GP with whole genome  
23 sequence data far more computationally expensive with millions of SNPs. However, it has been  
24 demonstrated that a multiple chain markov chain monte carlo methods for Bayesian GP results  
25 computationally cost effective and accurate predictions [26]. Similarly [27] found that updating

1 the right hand side of the Bayesian GP equations over multiple SNPs may reduce the need of  
2 memory allocations and computing time. However, it could be argued that the inflated  
3 accuracies were due to SNPs of GWAS results obtained from whole dataset. These results  
4 therefore need to be interpreted with caution. In order to correct for this bias, we performed the  
5 GWAS based on only training samples (Table 1), to select SNPs. As shown in Table 1 the  
6 accuracies were found to be lower compared with the results of the full dataset.

7 The present study was designed to determine the effect of prior distributions in Bayesian GP in  
8 terms of prediction accuracy under different experimental designs. The relevance of genetic  
9 architecture in conjugate to the prior distributions clearly supported by the unselected data. The  
10 most obvious finding emerge from this study is that preselection of SNPs referring to genetic  
11 architecture of the phenotype may lower the needs of computational load. Consistent with the  
12 literature, SNPs selection process should be exercised on training populations in order to avoid  
13 falsely inflated accuracies.

14

## LINKAGE DISEQUILIRIUM

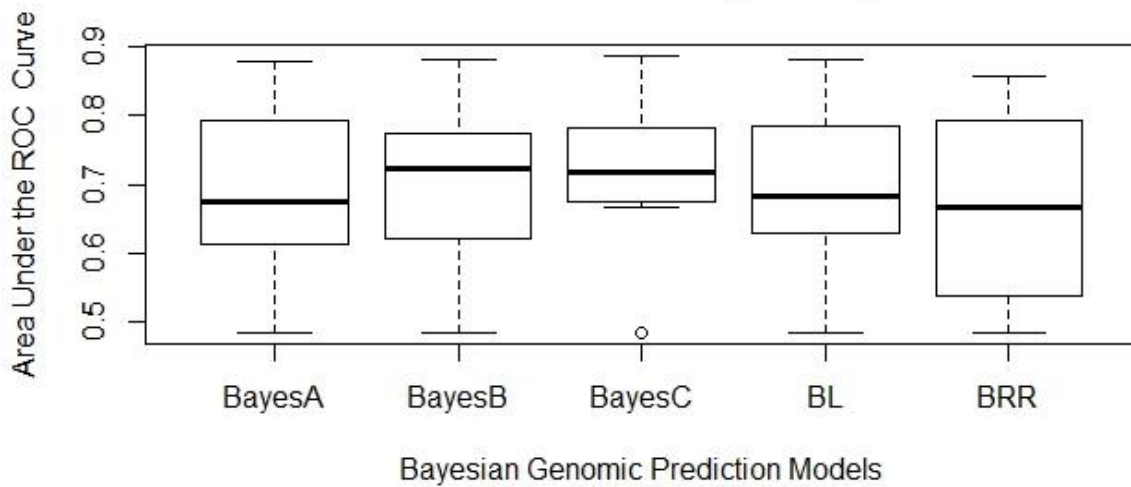


1

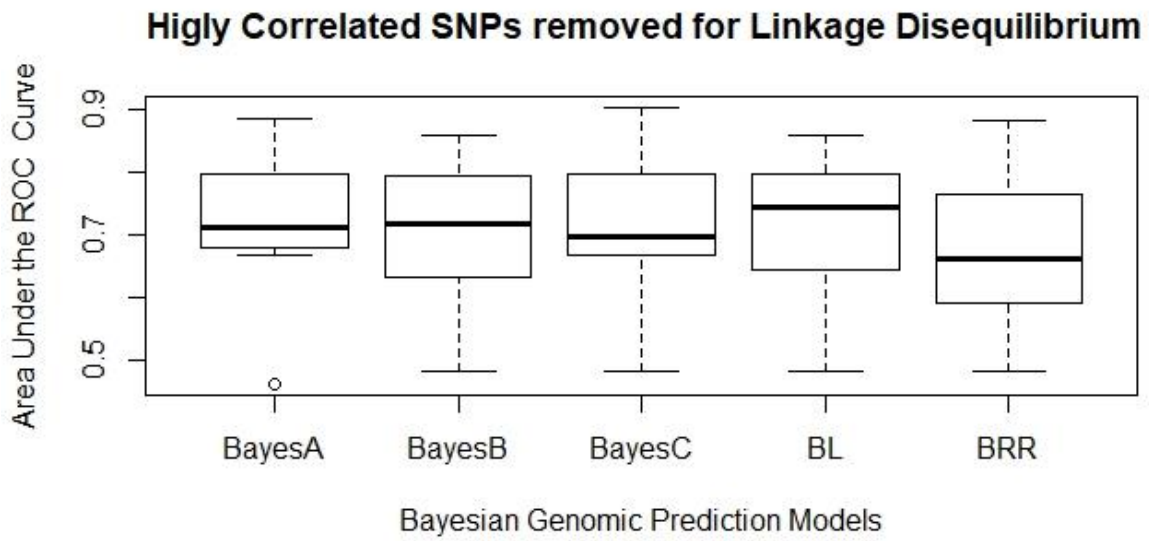
2 Figure 1. Decay of linkage disequilibrium over physical distance

3

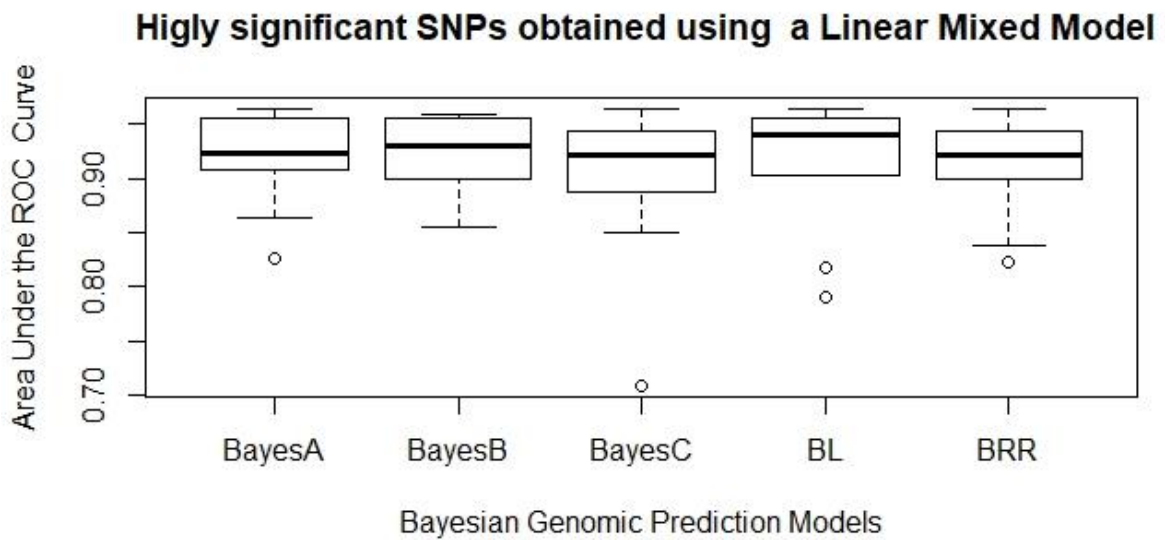
## No SNPs removed for Linkage Disequilibrium



4



1



2

3

4 Figure 2 Prediction accuracies obtained for Bayesian genomic prediction models from 10-fold  
 5 cross validations.

6



1 Table 1: Results of Bayesian learning models obtained from different experimental settings  
 2 over 10 fold cross validation procedure.

Model	No. SNPs	AUC (SD)
No SNPs removed for Linkage Disequilibrium		
BayesA	40642	0.684 (0.119)
BayesB	40642	0.710 (0.117)
BayesC	40642	0.724 (0.113)
BL	40642	0.700 (0.125)
BRR	40642	0.673 (0.130)
Highly Correlated SNPs removed for Linkage Disequilibrium ( $r^2 > 0.7$ )		
BayesA	25254	0.727 (0.118)
BayesB	25254	0.699 (0.132)
BayesC	25254	0.721 (0.119)
BL	25254	0.711 (0.134)
BRR	25254	0.679 (0.128)
SNPs selected by using Linear Mixed Model results ( $P < 0.05$ )		
BayesA	2401	0.919 (0.044)
BayesB	2401	0.922 (0.039)
BayesC	2401	0.898 (0.074)
BL	2401	0.915 (0.061)
BRR	2401	0.910 (0.047)
SNPs selected by using Linear Mixed Model results ( $P < 0.05$ ) obtained from training populations		
BayesA	2120	0.613 (0.126)
BayesB	2120	0.624 (0.124)
BayesC	2120	0.694 (0.081)
BL	2120	0.619 (0.114)
BRR	2120	0.673 (0.103)

3  
 4  
 5  
 6  
 7  
 8

## 1 References

- 2 1. VanRaden PM. Symposium review: How to implement genomic selection. *Journal of*  
3 *dairy science* 2020; 103.6: 5291-5301.
- 4 2. Gutierrez-Reinoso MA, Aponte PM, Garcia-Herreros M. Genomic Analysis, Progress  
5 and Future Perspectives in Dairy Cattle Selection: A Review. *Animals* 2021; 11.3: 599.
- 6 3. Meuwissen T, Hayes B, Goddard M. Genomic selection: A paradigm shift in animal  
7 breeding. *Animal frontiers* 2016; 6.1: 6-14.
- 8 4. Baker LA, Momen M, Chan K, Bollig N, Lopes FB, Rosa G J, Muir P. Bayesian and  
9 machine learning models for genomic prediction of anterior cruciate ligament rupture  
10 in the canine model. *G3: Genes, genomes, genetics* 2020; 10.8: 2619-2628.
- 11 5. Grinberg NF, Orhobor OI, King RD. An evaluation of machine-learning for predicting  
12 phenotype: studies in yeast, rice, and wheat. *Machine Learning* 2020; 109.2: 251-277.
- 13 6. Shi S, Li X, Fang L, Liu A, Su G, Zhang Y, Zhang S. Genomic Prediction Using  
14 Bayesian Regression Models With Global–Local Prior. *Frontiers in Genetics* 2021; 12:  
15 426.
- 16 7. Gianola D. Priors in whole-genome regression: the Bayesian alphabet returns. *Genetics*  
17 2013; 194.3: 573-596.
- 18 8. Suárez-Vega A, Gutiérrez-Gil B, Benavides J, Perez V, Tosser-Klopp G, Klopp C,  
19 Arranz JJ. Combining GWAS and RNA-Seq approaches for detection of the causal  
20 mutation for hereditary junctional epidermolysis bullosa in sheep. *PLoS One*  
21 2015; 10.5: e0126416.
- 22 9. Mömke S, Kerkmann A, Wöhlke A, Ostmeier M, Hewicker-Trautwein M, Ganter M,  
23 Distl O. A frameshift mutation within LAMC2 is responsible for Herlitz type junctional  
24 epidermolysis bullosa (HJEB) in black headed mutton sheep. *PloS one* 2011; 6.5:  
25 e18943.
- 26 10. Sartelet A, Harland C, Tamma N, Karim L, Bayrou C, Li W, Charlier C. A stop-gain  
27 in the laminin, alpha 3 gene causes recessive junctional epidermolysis bullosa in Belgian  
28 Blue cattle. *Animal genetics* 2015; 46.5: 566-570.
- 29 11. Kerkmann A, Ganter M, Frase R, Ostmeier M, Hewicker-Trautwein M, Distl O.  
30 Epidermolysis bullosa in German black headed mutton sheep. *Berliner und Munchener*  
31 *Tierärztliche Wochenschrift* 2010; 123.9-10: 413-421.
- 32 12. Ostmeier M, Kerkmann A, Frase R, Ganter M, Distl O, Hewicker-Trautwein M.  
33 Inherited junctional epidermolysis bullosa (Herlitz type) in German black-headed  
34 mutton sheep. *Journal of comparative pathology* 2012; 146.4: 338-347.

- 1 13. Milenkovic D, Chaffaux S, Taourit S, Guérin G. A mutation in the LAMC2 gene causes  
2 the Herlitz junctional epidermolysis bullosa (H-JEB) in two French draft horse  
3 breeds. *Genetics Selection Evolution* 2003; 35.2: 249-256.
- 4 14. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MA, Bender D, Sham PC.  
5 PLINK: a tool set for whole-genome association and population-based linkage  
6 analyses. *The American journal of human genetics* 2007; 81.3: 559-575.
- 7 15. Meuwissen TH, Hayes BJ, Goddard ME. Prediction of total genetic value using  
8 genome-wide dense marker maps. *Genetics* 2001; 157.4: 1819-1829.
- 9 16. Park T, Casella G. The bayesian lasso. *Journal of the American Statistical Association*  
10 2008; 103.482: 681-686.
- 11 17. Moser G, Lee SH, Hayes BJ, Goddard ME, Wray NR, Visscher PM. Simultaneous  
12 discovery, estimation and prediction analysis of complex traits using a Bayesian mixture  
13 model. *PLoS genetics* 2015; 11.4: e1004969.
- 14 18. Pérez P, de Los Campos G. Genome-wide regression and prediction with the BGLR  
15 statistical package. *Genetics* 2014; 198.2: 483-495.
- 16 19. Schulz-Streeck T, Ogotu JO, Piepho HP. Pre-selection of markers for genomic  
17 selection. In *BMC proceedings* 2011; 5.3: 1-4. BioMed Central.
- 18 20. Zhou X, Stephens M. Genome-wide efficient mixed-model analysis for association  
19 studies. *Nature genetics* 2012; 44.7: 821-824.
- 20 21. Calus MPL, De Haas Y, Veerkamp RF. Combining cow and bull reference populations  
21 to increase accuracy of genomic prediction and genome-wide association  
22 studies. *Journal of Dairy Science* 2013; 96.10: 6703-6715.
- 23 22. Fanny M, Andrea R, Pascal C. Evaluating the Interpretability of SNP Effect Size Classes  
24 in Bayesian Genomic Prediction Models. In *HUMAN HEREDITY* 2021; 85.2: 86-86.
- 25 23. Bang SJ, Kim YG, Park T. Joint selection of SNPs for improving prediction in genome-  
26 wide association studies. In *2012 IEEE International Conference on Bioinformatics and*  
27 *Biomedicine Workshops* 2012: 852-858.
- 28 24. Xiang R, Breen E, Prowse-Wilkins C, Chamberlain A, Goddard, M. Bayesian genome-  
29 wide analysis of cattle traits using variants with functional and evolutionary  
30 significance. *bioRxiv* 2021.
- 31 25. Meuwissen T, van den Berg I, Goddard M. On the use of whole-genome sequence data  
32 for across-breed genomic prediction and fine-scale mapping of QTL. *Genetics Selection*  
33 *Evolution* 2021; 53.1: 1-15.
- 34 26. Guo P, Zhu B, Niu H, Wang Z, Liang Y, Chen Y, Li J. Fast genomic prediction of  
35 breeding values using parallel Markov chain Monte Carlo with convergence  
36 diagnosis. *BMC bioinformatics* 2018; 19.1: 1-11.

- 1 27. Calus MP. Right-hand-side updating for fast computing of genomic breeding
- 2 values. *Genetics Selection Evolution* 2014; 46.1: 1-11.
- 3