

# 1 **Bayesian genome-wide analysis of cattle traits using variants with** 2 **functional and evolutionary significance**

3

4 Ruidong Xiang<sup>1,2,\*</sup>, Ed J. Breen<sup>2</sup>, Claire P. Prowse-Wilkins<sup>1,2</sup>, Amanda J. Chamberlain<sup>2</sup>,  
5 Michael E. Goddard<sup>1,2</sup>

6 <sup>1</sup> *Faculty of Veterinary & Agricultural Science, The University of Melbourne, Parkville 3052,*  
7 *Victoria, Australia*

8 <sup>2</sup> *Agriculture Victoria, AgriBio, Centre for AgriBiosciences, Bundoora, Victoria 3083, Australia.*

9 \* Corresponding author: [ruidong.xiang@unimelb.edu.au](mailto:ruidong.xiang@unimelb.edu.au)

10

## 11 **Abstract**

12 **Context.** Functional genomics studies have revealed genomic regions with regulatory and  
13 evolutionary significance. Such information independent of association analysis may benefit  
14 fine-mapping and genomic selection of economically important traits. However, systematic  
15 evaluation of the use of functional information in mapping, and genomic selection of cattle  
16 traits is lacking. Also, Single Nucleotide Polymorphisms (SNPs) from the high-density (HD)  
17 panel are known to tag informative variants, but the performance of genomic prediction using  
18 HD SNPs together with variants supported by different functional genomics is unknown.

19 **Aims.** We selected six sets of functionally important variants and modelled each set together  
20 with HD SNPs in Bayesian models to map and predict protein, fat, and milk yield as well as  
21 mastitis, somatic cell count and temperament of dairy cattle.

22 **Methods.** Two models were used: 1) BayesR which includes priors of four distribution of  
23 variant-effects, and 2) BayesRC which includes additional priors of different functional  
24 classes of variants. Bayesian models were trained in 3 breeds of 28,000 cows of Holstein,  
25 Jersey and Australian Red and predicted into 2,600 independent bulls.

26 **Key results.** Adding functionally important variants significantly increased the enrichment of  
27 genetic variance explained for mapped variants, suggesting improved genome-wide mapping  
28 precision. Such improvement was significantly higher when the same set of variants were  
29 modelled by BayesRC than by BayesR. Combining functional variant sets with HD SNPs  
30 improves genomic prediction accuracy in the majority of the cases and such improvement  
31 was more common and stronger for non-Holstein breeds and traits like mastitis, somatic cell  
32 count and temperament. In contrast, adding a large number of random sequence variants to  
33 HD SNPs reduces mapping precision and has a worse or similar prediction accuracy,  
34 compared to using HD SNPs alone to map or predict. While BayesRC tended to have better  
35 genomic prediction accuracy than BayesR, the overall difference in prediction accuracy  
36 between the two models was insignificant.

37 **Conclusions.** Our findings demonstrate the usefulness of functional data in genomic mapping  
38 and prediction.

39 **Implications.** We highlight the need for effective tools exploiting complex functional  
40 datasets to improve genomic prediction.

41

42 **Key words:** Functional genomics, Animal breeding, Genetic mapping, Quantitative genetics

43

## 44 **Introduction**

45 Emerging evidence shows that genomic variants with causal roles in biology can be used to  
46 improve genomic prediction of complex traits. The biological function of genomic variants  
47 provides information independent of genotype-trait associations which are usually  
48 confounded by linkage disequilibrium (LD). Such independent information can be exploited  
49 to identify informative variants. Once identified, informative variants can be used to improve  
50 genomic prediction (Xiang, MacLeod, Daetwyler, de Jong, O'Connor, Schrooten,  
51 Chamberlain & Goddard, 2021). While the use of functional data in improving genomic  
52 mapping and prediction has been reported in humans (Amariuta, Ishigaki, Sugishita, Ohta,  
53 Koido, Dey, Matsuda, Murakami, Price & Kawakami, 2020; Weissbrod, Hormozdiari,  
54 Benner, Cui, Ulirsch, Gazal, Schoech, Van De Geijn, Reshef & Márquez-Luna, 2020), using  
55 functional data in predicting the genetic merit of animal traits has not been comprehensively  
56 examined. However, there is evidence in cattle supporting the advantage of the use of  
57 functional information in genomic mapping and prediction with the linear mixed model  
58 (Fang, Sahana, Ma, Su, Yu, Zhang, Lund & Sørensen, 2017a; Fang, Sahana, Ma, Su, Yu,  
59 Zhang, Lund & Sørensen, 2017b; Liu, Fang, Zhou, Santos, Xiang, Daetwyler, Chamberlain,  
60 Cole, Li, Yu, Ma, Zhang & Liu, 2019; Xiang, Berg, MacLeod, Hayes, Prowse-Wilkins,  
61 Wang, Bolormaa, Liu, Rochfort, Reich, Mason, Vander Jagt, Daetwyler, Lund, Chamberlain  
62 & Goddard, 2019; Xu, Gao, Wang, Xu, Liu, Chen, Xu, Gao, Zhang & Gao, 2020).

63 The Functional Annotation of ANimal Genomes (FAANG) consortium (Clark, Archibald,  
64 Daetwyler, Groenen, Harrison, Houston, Kühn, Lien, Macqueen & Reecy, 2020) provides  
65 many types of sequencing data indicating the functionality of genome-wide sites (examples  
66 reviewed in (Clark *et al.*, 2020)). While these public datasets await exploitation, the structure  
67 and information content of different functional datasets vary significantly. For example, we  
68 recently showed that amongst all analysed functional datasets, a set of 300,000+ sequence

69 variants within sites highly conserved across 100 vertebrate species had the strongest  
70 enrichment with cattle trait heritability (Xiang *et al.*, 2019), which primarily influences  
71 genomic prediction accuracy. Additionally, a few thousand variants affecting the  
72 concentration of milk fat metabolites, i.e., metabolic mQTLs, also had significantly higher  
73 variance than SNPs in the 50K panel for cattle traits. Millions of variants that change gene  
74 expression levels (geQTLs) or RNA splicing (sQTLs) are also enriched with complex trait  
75 QTL (Fink, Lopdell, Tiplady, Handley, Johnson, Spelman, Davis, Snell & Littlejohn, 2020;  
76 Li, van de Geijn, Raj, Knowles, Petti, Golan, Gilad & Pritchard, 2016; Lopdell, Tiplady,  
77 Struchalin, Johnson, Keehan, Sherlock, Couldrey, Davis, Snell & Spelman, 2017; Silva,  
78 Fonseca, Pinheiro, Magalhães, Muniz, Ferro, Baldi, Chardulo, Schnabel & Taylor, 2020;  
79 Xiang, Hayes, Vander Jagt, MacLeod, Khansefid, Bowman, Yuan, Prowse-Wilkins, Reich,  
80 Mason, Garner, Marett, Chen, Bolormaa, Daetwyler, Chamberlain & Goddard, 2018).  
81 However, recent studies showed that variants close to genes with high or specific expression  
82 patterns had limited improvement in prediction accuracy (de Las Heras-Saldana, Lopez,  
83 Moghaddar, Park, Park, Chung, Lim, Lee, Shin & van der Werf, 2020; Fang, Cai, Liu,  
84 Canela-Xandri, Gao, Jiang, Rawlik, Li, Schroeder & Rosen, 2020). Another common type of  
85 functional data is peaks from ChIP-seq for histone modifications which are enriched with  
86 promoters and/or enhancers regulating gene activities (Carey, Peterson & Smale, 2009). Our  
87 work showed that hundreds of thousands of variants under ChIP-seq peaks are enriched for  
88 complex trait QTL in cattle (Prowse-Wilkins, Wang, Xiang, Goddard & Chamberlain, 2021;  
89 Xiang *et al.*, 2019). In addition, variants within the gene coding regions are expected to have  
90 a high impact on complex traits. However, we and others previously found coding-related  
91 variants (around 100,000) have limited contributions to cattle trait heritability (Koufariotis,  
92 Chen, Stothard & Hayes, 2018; Xiang *et al.*, 2019), although their use in improving genomic  
93 prediction has not been studied.

94 One way to assess the information content of functional data is to compare variants  
95 prioritised by functional data with SNPs from standard genotyping panels. We have  
96 previously performed such assessment using the standard 50K bovine SNP chip and showed  
97 that functional information can improve genomic prediction accuracy compared to the 50K  
98 chip SNPs (Xiang *et al.*, 2021). However, denser panels such as the high-density (HD) SNP  
99 chip containing ~700,000 SNPs across the genome may be able to tag many functional  
100 elements via LD, although it is not routinely used in animal genomic evaluation. With the  
101 development of animal breeding, the HD panel may be intensively used in the future genomic  
102 evaluation. Therefore, it is of interest to know if functional information can provide any  
103 advantage in genomic mapping and prediction when HD SNPs are used. Also, since causal  
104 variants are expected to have similar phenotypic effects across different breeds, we aim to  
105 compare the use of functionally important variants in genomic prediction across different  
106 breeds.

107 In the present study, we evaluate sequence variant sets prioritised by 6 types of functional and  
108 evolutionary data in combination with the standard HD SNPs in genomic mapping and  
109 prediction of 6 dairy cattle traits. We train the prediction equations using the BayesR method  
110 (Erbe, Hayes, Matukumalli, Goswami, Bowman, Reich, Mason & Goddard, 2012) which fits  
111 a mixture of 4 distributions of variant-effects and using the BayesRC method which fits  
112 different distributions for each functional class of variant classifications (MacLeod, Bowman,  
113 Vander Jagt, Haile-Mariam, Kemper, Chamberlain, Schrooten, Hayes & Goddard, 2016).

114 Genomic predictors were trained using 28,000 cows that included 3 breeds: Holstein, Jersey  
115 and Australian Red. Genomic estimated breeding values (gEBVs) were predicted and  
116 validated in 2,500 Holstein, Jersey and Australian Red bulls. We compare the results of  
117 mapping and genomic prediction across the above-described scenarios, discuss these results  
118 and provide suggestions for future studies.

119

## 120 **Materials and Methods**

121 The phenotype data analysed in this study were collected by DataGene Australia  
122 (<http://www.datagene.com.au/>) and no further live animal experimentation was required for  
123 our analyses. A set of 28,049 Australian cows were used as the discovery population and a set  
124 of 2,567 bulls were used as the validation population. The bull phenotypes were obtained as  
125 daughter trait deviations: i.e. the average trait deviations of a bull's daughters pre-corrected  
126 for known fixed effects by DataGene. The cow phenotypes were measured on themselves.  
127 Note that these bulls and cows were not included in those 44,000+ animals used to discover  
128 functional variants (Xiang *et al.*, 2019; Xiang *et al.*, 2021; Xiang, van den Berg, MacLeod,  
129 Daetwyler & Goddard, 2020). We also checked the pedigree to make sure that bulls used in  
130 the validation population were not the sires of cows from the discovery population. Cows in  
131 the discovery set included 24,305 Holstein, 2,486 Jersey, 1,258 Australian Red. Bulls in the  
132 validation datasets contained 2,091 Holstein, 385 Jersey, 91 Australian Red. Traits  
133 considered in the analysis included protein yield (Prot), fat yield (Fat), milk yield (Milk),  
134 Mastitis (Mas), somatic cell count (Scc) and temperament (Temp).  
135 The genotypes used in the study were imputed sequence variants based on Run7 of the 1000  
136 Bull Genomes Project (Daetwyler, Capitan, Pausch, Stothard, Van Binsbergen, Brøndum,  
137 Liao, Djari, Rodriguez & Grohs, 2014; Hayes & Daetwyler, 2018) based on the ARS-  
138 UCD1.2 reference bovine genome  
139 ([https://www.ncbi.nlm.nih.gov/assembly/GCF\\_002263795.1/](https://www.ncbi.nlm.nih.gov/assembly/GCF_002263795.1/)) (Rosen, Bickhart, Schnabel,  
140 Koren, Elsik, Tseng, Rowan, Low, Zimin & Couldrey, 2020). Variants with Minimac3  
141 (Fuchsberger, Abecasis & Hinds, 2014; Howie, Fuchsberger, Stephens, Marchini & Abecasis,  
142 2012) imputation accuracy  $R^2 > 0.4$  and minor allele frequency (MAF)  $> 0.005$  in bulls and  
143 cows. Most bulls were genotyped with a medium-density SNP array (50K) or a high-density

144 SNP array and most cows were genotyped with a low-density panel of approximately 6,900  
145 SNPs overlapping with the standard-50K panel (BovineSNP50 beadchip, Illumina Inc). The  
146 low-density genotypes were first imputed to the Standard-50K panel and then all 50K  
147 genotypes were imputed to the HD panel using Fimpute v3 (Sargolzaei, Chesnais &  
148 Schenkel, 2014; Xiang *et al.*, 2019). Then, all HD genotypes were imputed to sequence using  
149 Minimac3 with Eagle (v2) to pre-phase genotypes (Howie *et al.*, 2012; Loh, Danecek,  
150 Palamara, Fuchsberger, Reshef, Finucane, Schoenherr, Forer, McCarthy & Abecasis, 2016).  
151 We aimed to test whether variant sets selected from different functional and/or evolutionary  
152 information, in addition to the standard HD SNP panel, can be useful for genomic prediction.  
153 Therefore, we first included a baseline set, which is 610,764 SNPs from the standard bovine  
154 high-density panel. There were six functional and/or evolutionary variant sets: 549,007  
155 variants under multiple ChIP-seq peaks (Kern, Wang, Xu, Pan, Halstead, Chanthavixay,  
156 Saelao, Waters, Xiang & Chamberlain, 2021; Prowse-Wilkins *et al.*, 2021) ('ChiPseq'),  
157 106,538 variants annotated as related to coding activities by Ensembl Variant Effect Predictor  
158 (McLaren, Gil, Hunt, Riat, Ritchie, Thormann, Flicek & Cunningham, 2016) ('Coding'),  
159 943,315 variants affecting RNA splicing sQTLs from 4 cattle tissues (Chamberlain, Hayes,  
160 Xiang, Vander Jagt, Reich, Macleod, Prowse-Wilkins, Mason, Daetwyler & Goddard, 2018;  
161 Daetwyler, Xiang, Yuan, Bolormaa, Vander Jagt, Hayes, van der Werf, Pryce, Chamberlain  
162 & Macleod, 2019; Xiang *et al.*, 2018) ('sQTL'), 65,394 finely mapped variants with  
163 pleiotropic effects genome-wide (Xiang *et al.*, 2021) ('Finemap80k'), 4,871 variants affecting  
164 milk fat metabolites mQTLs (Xiang *et al.*, 2019) ('mQTL') and 317,279 conserved sites  
165 across 100 vertebrates (Xiang *et al.*, 2019) ('Cons100w'). Note that some of these functional  
166 variant sets were initially determined on the UMD3.1 genome and were from different cattle  
167 populations. These sets were lifted over from the older genome to ARS-UCD1.2 and filtered  
168 with imputation accuracy and MAF in the new cattle populations.

169 The model training of the above-described data used BayesR (Erbe *et al.*, 2012) and  
170 BayesRC (MacLeod *et al.*, 2016), which are now implemented via BayesR3, with improved  
171 efficiency using blocks. BayesR jointly models all variants together with different effect  
172 distribution priors. BayesRC follows the same approach but in addition allows a ‘C’ prior  
173 which models classes of variants. Another aim is to see whether there are differences in  
174 genomic prediction accuracy by modelling the same variants using BayesR and BayesRC. To  
175 aid this comparison, we combined each functional variant set with the HD variants which led  
176 to 6 combined variant sets: 1) ChIP-seq peak tagged variants + HD SNPs (‘ChiPseq\_HD’), 2)  
177 coding variants + HD variants (‘Coding\_HD’), 3) sQTL variants + HD SNPs (‘sQTL\_HD’),  
178 4) finely mapped variants + HD SNPs (‘Finemap80k\_HD’), 5) mQTL variants + HD SNPs  
179 (‘mQTL\_HD’) and 6) conserved variants + HD SNPs (‘Cons100w\_HD’). The average minor  
180 allele frequency of these sets of variants were 0.22 ( $\pm 0.00014$ ) for ChiPseq\_HD,  
181 0.25 ( $\pm 0.0002$ ) for Coding\_HD, 0.24 ( $\pm 0.0001$ ) for sQTL\_HD, 0.27 ( $\pm 0.0002$ ) for  
182 Finemap80k\_HD, 0.27 ( $\pm 0.0002$ ) for mQTL\_HD, 0.23 ( $\pm 0.0002$ ) for Cons100w\_HD, and  
183 0.27 ( $\pm 0.0002$ ) for HD alone.

184 In single-trait BayesR, we directly model these 6 variant sets one set at a time. To create a  
185 reference baseline, we also used single-trait BayesR to fit the HD variant set (‘HD’) alone. In  
186 single-trait BayesRC, for each of the same 6 combined variant sets, we specify 2 different  
187 variant classes: 1) Variants appeared in the functional and/or evolutionary set and 2) variants  
188 only appeared in the HD variant set.

189 Both BayesR and BayesRC modelled variant effects as a mixture distribution of four normal  
190 distributions including a null distribution,  $N(0, 0.0\sigma_g^2)$ , and three others:  $N(0, 0.0001\sigma_g^2)$ ,  
191  $N(0, 0.001\sigma_g^2)$ ,  $N(0, 0.01\sigma_g^2)$ , where  $\sigma_g^2$  was the additive genetic variance for the trait.

192 The starting value of  $\sigma_g^2$  for each trait was estimated using GREML implemented in the



193 MTG2 (Lee & Van der Werf, 2016) with a single genomic relationship matrix made of all  
194 sequence variants. The statistical model used in the single-trait BayesR and BayesRC in was:

195 
$$\mathbf{y} = \mathbf{W}\mathbf{v} + \mathbf{X}\mathbf{b} + \mathbf{e} \text{ (equation 1)}$$

196 where  $\mathbf{y}$  was a vector of phenotypic records;  $\mathbf{W}$  was the design matrix of marker genotypes;  
197 centred and standardised to have a unit variance;  $\mathbf{v}$  was the vector of variant effects,  
198 distributed as a mixture of the four distributions as described above;  $\mathbf{X}$  was the design matrix  
199 allocating phenotypes to fixed effects;  $\mathbf{b}$  was the vector of fixed effects, including breeds;  
200  $\mathbf{e}$  = vector of residual errors. As a result, the effect  $b$  for each variant jointly estimated with  
201 other variants were obtained for further analysis.

202 BayesRC used the same linear model as BayesR. The C component of BayesRC had two  
203 categories  $c$  ( $c = 2$ ) as described above. Within each category  $c$ , an uninformative Dirichlet  
204 prior ( $\alpha$ ) was used for the proportion of effects in each of the four normal distributions of  
205 variant effects:  $P_c \sim Dir(\alpha_c)$ , where  $\alpha_c = [1, 1, 1, 1]$ .  $\alpha_c$  was updated each iteration within  
206 each category:  $P_c \sim Dir(\alpha_c + \beta_c)$ , where  $\beta_c$  was the current number of variants in each of the  
207 four distributions within category  $c$ , as estimated from the data.

208 Two metrics were evaluated for mapping results. One is the mixing proportion, i.e., the  
209 proportion of variants with small effect  $N(0, 0.0001\sigma_g^2)$ , medium effect  $N(0, 0.001\sigma_g^2)$   
210 and large effect  $N(0, 0.01\sigma_g^2)$  for each BayesRC run across the functional variant class and  
211 the HD SNP class. This metric shows the information content of the two classes. The other  
212 metric was the percentage of 50kb segments needed by the model to explain 50% of the  
213 cumulative sum of posterior probability (PP), which indicated the mapping precision. For  
214 each variant, PP was calculated as  $1 - P_0$  where  $P_0$  was the probability for the variant to be  
215 within the zero-effect distribution  $N(0, 0.0\sigma_g^2)$ . The sum of PP across all variants estimates  
216 the number of variants causing genetic variance in the trait. The smaller amount of genomic  
217 segments needed to explain a cumulative sum of PP, the higher the mapping precision. We

218 also compared genomic prediction accuracy, defined as the Pearson correlation  $r$  between  
219 genomic estimated breeding value (gEBV) and phenotype in the validation populations.  
220 gEBV of the validation animals was calculated as  $gEBV = \mathbf{Z}\hat{\mathbf{s}}$  (equation 2), where  $\mathbf{Z}$  was a  
221 matrix of the standardised genotypes of animals in the validation set, and  $\hat{\mathbf{s}}$  was the vector of  
222 variant effects from the training model. In addition, to test if adding a large number of  
223 random variants to the HD panel can increase mapping precision and prediction accuracy, a  
224 random set of 944,616 variants matching the size of the largest set of functional variants  
225 (sQTL, 943,315 variants) was also selected and added to the HD panel ('Random\_HD'). This  
226 random set was analysed for BayesR, mapping precision and prediction accuracy in the same  
227 fashion as other variant sets described above.

228

## 229 **Results**

### 230 *Information content in the functional variant sets*

231 Averaged across mixing proportions from single-trait BayesRC, we show that compared to  
232 HD SNPs, the finely mapped variants had consistently higher enrichment with variants  
233 showing small, medium and large effects (Figure 1). Variants within coding regions showed  
234 higher enrichment than HD SNPs for large- and medium-effect variants. Interestingly,  
235 mQTLs, which were variants affecting the concentration of milk fat metabolites (Benedet,  
236 Ho, Xiang, Bolormaa, De Marchi, Goddard & Pryce, 2019; Xiang *et al.*, 2019), had lower  
237 enrichment of small-effect variants than HD SNPs, but had higher enrichment of medium and  
238 large-effect variants than HD SNPs.

### 239 *Mapping precision*

240 Across traits, we show that all models using functional variants, except mQTL, needed a  
241 smaller amount of genome-wide segments to explain 50% of the cumulative sum of PP,  
242 compared to HD SNPs (Figure 2). This means that when adding to the HD SNPs, most

243 functional variants increased mapping precision. In contrast, adding randomly selected  
244 944,000 variants to HD SNPs increased the amount of genome-wide segments (by  $2.82\% \pm$   
245  $0.13\%$ ) across scenarios to explain 50% of the cumulative sum of PP, compared to only using  
246 HD SNPs. This suggested that adding random variants to HD decreases mapping precision. It  
247 is worth noting that when using 106,538 coding variants and 65,394 finely mapped variants,  
248 BayesRC provided a further increase in mapping precision over HD SNPs than BayesR. On  
249 the other hand, when using 549,007 ChIP-seq tagged variants and 943,315 sQTL variants,  
250 BayesRC had less increase in mapping precision over HD SNPs than BayesR. This could be  
251 due to the reduced signal-to-noise ratio in large variant sets of ChIP-seq tagged variants and  
252 sQTLs.

### 253 *Genomic prediction of traits*

254 In total, we evaluated the genomic prediction accuracy in 216 scenarios, across 6 single-trait  
255 analysis, 6 functional categories, 4 breeds in the validation population, and 2 Bayesian  
256 methods. Out of these 216 scenarios, 142 (66%) times, HD SNPs combined with functional  
257 variants increased genomic prediction accuracy, compared to the prediction only using the  
258 HD SNPs (Figure 3 and 4). In 51 out of 216 times (24%), the increase in prediction accuracy  
259 ( $[r_{functional} - r_{HD}] \times 100\%$ ) was greater than 1%. These 51 cases were almost all accounted  
260 for by Jersey (15/51) and Australian Red (34/51), with only 2 cases in Holstein cattle. In 29  
261 analyses (14%), the increase in prediction accuracy over HD SNPs was greater than 2%. All  
262 these 29 cases were for non-Holstein breeds. Amongst tested functional sets, genomic  
263 prediction accuracy was the best when the HD variants were combined with conserved  
264 variants (Cons100w\_HD). In contrast, averaged across tested scenarios, adding randomly  
265 selected 944,000 variants to HD had a slightly worse or no improvement in prediction  
266 accuracy ( $-0.5\% \pm 0.49\%$ ) compared to only using the HD panel to predict.

267 As shown in Figure 3, the genomic prediction accuracy of milk production traits using HD  
268 SNPs in Holstein cattle was already high (around 0.7) and the increases in accuracy from  
269 functional variants were very small. However, larger increases were evident in Jersey and  
270 Australian Red. For milk production traits, 10 out of 18 times the genomic prediction  
271 accuracy was the most improved by conserved variants and coding variants combined with  
272 HD SNPs, followed by finely mapped variants combined with HD SNPs (4/18), ChIP-seq  
273 tagged variants (3/18) combined with HD SNPs. sQTL combined with HD variants had the  
274 highest accuracy when predicting protein yield in Holstein.

275 As shown in Figure 4, the greatest increases in prediction accuracy for traits mastitis, somatic  
276 cell count and temperament were again seen in non-Holstein breeds. Chip-seq peak tagged  
277 variants combined with HD SNPs (5/18 times) and conserved variants combined with HD  
278 SNPs (5/18 times) had the best performances in predicting mastitis, somatic cell count and  
279 temperament.

280 Across all scenarios, we did not see a clear distinction in prediction accuracy between  
281 BayesR and BayesRC in the current study. There may be some tendencies where BayesRC  
282 had a higher accuracy than BayesR for somatic cell count, mastitis and temperament.  
283 However, none of these differences were significant.

284

## 285 **Discussion**

286 Our systematic evaluations show that functional information can improve genomic mapping  
287 and prediction of cattle traits, even when HD SNPs are used, although there were times where  
288 HD SNPs alone still had robust performances. It is usually the less represented breeds, such  
289 as Jersey and Australian Red who benefited the most from the improvements using functional  
290 data. This suggests functional information can well complement HD SNPs especially in  
291 breeds with smaller training sets. Adding randomly selected variants to the HD panel reduces

292 mapping precision and provided no improvement in prediction accuracy compared to only  
293 using the HD panel. This supports that the The benefit provided via selecting variants based  
294 on functional importance can not simply be achieved by adding more sequence variants.  
295 We show that the biological information content which can be used to benefit mapping and/or  
296 prediction is different between functional datasets. One of the top-performing functional  
297 variant sets in mapping large-effect variants was the finely mapped 80,000 variants (Xiang *et*  
298 *al.*, 2021). This result is somewhat expected as these variants combined information from  
299 multiple functional datasets and also included variants affecting multiple dairy cattle traits.  
300 These finely mapped 80,000 variants outperformed the SNPs from the 50K panel in previous  
301 evaluations (Xiang *et al.*, 2021). Furthermore, finely mapped 80,000 variants showed  
302 enhanced enrichment of large-effect variants and improvement in mapping precision when  
303 modelled with BayesRC. This suggests that this much more refined set of variants (chosen  
304 because they were more relevant to the traits of interest) are likely more enriched for variants  
305 that are more strongly associated with the trait or are causal. BayesRC would only  
306 outperform BayesR when there is strong enrichment for QTL in at least one of the defined  
307 classes. The other functional groups tested are not trait specific (except mQTL for fat) so  
308 likely less enriched relative to each trait.

309 Previous results showed that coding-related variants did not explain a significant amount of  
310 heritability (Koufariotis *et al.*, 2018; Xiang *et al.*, 2019). In the current study, coding-related  
311 variants combined with HD SNPs showed enhanced enrichment with large-effect variants  
312 and improvement in mapping precision. This implies that variants affecting protein coding  
313 may not necessarily be good at capturing all the genetic variance of polygenic traits. The  
314 small set of mQTLs, derived from milk fat showed strong enrichment of large-effect variants  
315 but did not show improvement in mapping precision over HD SNPs. This set of variants  
316 needs future investigations.

317 Unlike the results in mapping large-effect variants, for genomic prediction, the top-  
318 performing variant set is the conserved variants combined with HD SNPs. The advantage of  
319 adding conserved variants to HD SNPs was particularly evident when predicting somatic cell  
320 count, mastitis and temperament of non-Holstein breeds (Figure 4). In fact, in these scenarios  
321 HD SNPs alone did not perform so well and this leaves more room for functional variants to  
322 improve the prediction accuracy. Another variant set that performed well in genomic  
323 prediction is the set of ChIP-seq peak tagged variants. Again, such an advantage was the most  
324 evident when predicting somatic cell count, mastitis and temperament in non-Holstein breeds.  
325 Interestingly, ChIP-seq variants combined with HD SNPs appear to show some particular  
326 advantages in predicting temperament. There may be some large-effect variants for  
327 temperament captured by ChIP-seq peaks.

328 We found that sQTL variants combined with HD SNPs had variable performances in  
329 mapping and prediction. This set did not show good performance in detecting enrichment of  
330 informative variants, but overall significantly increased mapping precision over HD SNPs. In  
331 genomic prediction, its performance was not impressive. This is somewhat different from  
332 previous studies which showed that sQTLs are enriched with complex trait QTL (Li *et al.*,  
333 2016; Xiang *et al.*, 2019; Xiang *et al.*, 2018). One explanation is that sQTLs or any other  
334 eQTLs were not trait specific and are plagued by LD, which is particularly strong for  
335 Holstein breeds that dominated the discovery population. Another explanation is that the  
336 sample size with which we used to discover sQTLs is still small (N~120) and we should re-  
337 discover and re-evaluate this set of variants when there is a larger sample size.

338 As mentioned earlier, BayesRC would only outperform BayesR when there is strong  
339 enrichment for QTL in at least one of the defined classes. It would also require functional  
340 information to be trait-specific. We saw advantages in BayesRC over BayesR in detecting  
341 enrichment with large-effect variants using finely mapped variants, coding variants and

342 mQTLs. BayesRC also had advantages over BayesR in mapping precision when used with  
343 finely mapped variants and coding variants. While these functional data are expected to be  
344 informative, they did not provide consistent advantages for BayesRC to predict traits over  
345 BayesR. Across all tested cases, we did not see strong advantages in BayesRC over BayesR  
346 in genomic prediction (Figure 4). BayesRC may have some tendencies to better predict  
347 somatic cell count, mastitis and temperament than BayesR. However, the differences were  
348 not statistically significant. The reason behind these observations may be complex.

349 We know that not all variants in the functional datasets are informative and many sequence  
350 variants are in strong LD. BayesR and BayesRC both have limitations where variants are in  
351 very strong LD. In addition, if most causal variants are quite well tagged by HD variants and  
352 if validation animals are highly related to the discovery animals, the room to improve  
353 prediction accuracy is limited. Also, there may be less common variants that are not tagged  
354 by HD SNPs, but these variants are not well imputed. Further, the optimal tissues and/or  
355 experimental conditions to generate functional data that can be better used for improving  
356 genomic prediction are usually not known. Therefore, the marriage between functional data  
357 and genomic prediction is still at its very early stage.

358 We therefore suggest two future research directions to improve on the current results. The  
359 first is to increase the information content in functional datasets. This can be achieved by  
360 either increasing the sample size (biological replicates, tissues and experimental conditions)  
361 of functional datasets or by developing better bioinformatic tools to increase the signal-to-  
362 noise ratio in functional datasets before they can be processed by genomic prediction models.

363 The second direction is to improve the current genomic prediction models. Because the type  
364 and complexity of functional data will keep growing, it will be necessary to develop more  
365 sophisticated and flexible methods to better extract information from complex functional  
366 data. For example, an extended BayesRC that can model quantitative biological priors,

367 instead of qualitative classes will be needed. Similarly, in the future we will use larger sample  
368 sizes and diverse breeds in the training model to reduce LD between sequence variants. This  
369 will also increase the need for Bayesian methods to be more efficient.

370 In conclusion, our evaluation of Bayesian genomic prediction using functional and  
371 evolutionary information with HD SNPs provides novel insights into this emerging area. We  
372 show that functional datasets of conserved variants, coding variants, ChIP-seq peaks and  
373 previously finely mapped variants can improve genomic mapping and/or genomic prediction,  
374 even when HD SNPs are used. Such improvements usually benefit non-Holstein breeds,  
375 given the current available functional datasets. We found that by using informative biological  
376 priors, BayesRC has significant advantages over BayesR in detecting enrichment with large-  
377 effect variants and in mapping precision. However, the advantage of BayesRC over BayesR  
378 for genomic prediction was not consistent. Our results highlight the need to develop better  
379 tools to extract information from complex functional datasets which will benefit genomic  
380 prediction in large datasets. Fusing functional genomics with genomic selection presents  
381 great opportunities to develop new technologies that improve animal breeding and genetics.

382

### 383 **Acknowledgments**

384 Australian Research Council's Discovery Projects (DP160101056 and DP200100499)  
385 supported R.X. and M.E.G. DairyBio, a joint venture project between Agriculture Victoria  
386 (Melbourne, Australia), Dairy Australia (Melbourne, Australia) and the Gardiner Foundation  
387 (Melbourne, Australia), funded computing resources used in the analysis. The authors also  
388 thank the University of Melbourne, Australia for supporting this research. No funding bodies  
389 participated in the design of the study nor analysis, or interpretation of data nor in writing the  
390 manuscript. DataGene and CRV provided access to the reference data used in this study. We  
391 thank Gert Nieuwhof, Kon Konstantinov and Timothy P. Hancock (DataGene) and staff from



392 DairyNZ for the preparation and provision of data. We thank Dr. Sunduimijid Bolormaa for  
393 the sequence variant data imputation. We thank Drs. Iona M. MacLeod and Hans D.  
394 Daetwyler for critical reading of the manuscript.

395

## 396 **Conflict of interest**

397 The authors declare no conflicts of interest.

398

## 399 **References:**

- 400 Amariuta, T., Ishigaki, K., Sugishita, H., Ohta, T., Koido, M., Dey, K.K., Matsuda, K.,  
401 Murakami, Y., Price, A.L. & Kawakami, E. 2020. Improving the trans-ancestry portability of  
402 polygenic risk scores by prioritizing variants in predicted cell-type-specific regulatory  
403 elements. *Nature Genetics*, 52(12):1346-1354.
- 404 Benedet, A., Ho, P., Xiang, R., Bolormaa, S., De Marchi, M., Goddard, M. & Pryce, J. 2019.  
405 The use of mid-infrared spectra to map genes affecting milk composition. *Journal of dairy  
406 science*, 102(8):7189-7203.
- 407 Carey, M.F., Peterson, C.L. & Smale, S.T. 2009. Chromatin immunoprecipitation (chip).  
408 *Cold Spring Harbor Protocols*, 2009(9):pdb. prot5279.
- 409 Chamberlain, A., Hayes, B., Xiang, R., Vander Jagt, C., Reich, C., Macleod, I., Prowse-  
410 Wilkins, C., Mason, B., Daetwyler, H. & Goddard, M. 2018. Identification of regulatory  
411 variation in dairy cattle with RNA sequence data.254.
- 412 Clark, E.L., Archibald, A.L., Daetwyler, H.D., Groenen, M.A., Harrison, P.W., Houston,  
413 R.D., Kühn, C., Lien, S., Macqueen, D.J. & Reecy, J.M. 2020. From FAANG to fork:  
414 application of highly annotated genomes to improve farmed animal production. *Genome  
415 Biology*, 21(1):1-9.
- 416 Daetwyler, H., Xiang, R., Yuan, Z., Bolormaa, S., Vander Jagt, C., Hayes, B., van der Werf,  
417 J., Pryce, J., Chamberlain, A. & Macleod, I. 2019. Integration of functional genomics and  
418 phenomics into genomic prediction raises its accuracy in sheep and dairy cattle. *Proceedings  
419 of the Association for the Advancement of Animal Breeding and Genetics, Armidale, NSW,  
420 Australia*:11-14.
- 421 Daetwyler, H.D., Capitan, A., Pausch, H., Stothard, P., Van Binsbergen, R., Brøndum, R.F.,  
422 Liao, X., Djari, A., Rodriguez, S.C. & Grohs, C. 2014. Whole-genome sequencing of 234  
423 bulls facilitates mapping of monogenic and complex traits in cattle. *Nature genetics*,  
424 46(8):858.
- 425 de Las Heras-Saldana, S., Lopez, B.I., Moghaddar, N., Park, W., Park, J.-e., Chung, K.Y.,  
426 Lim, D., Lee, S.H., Shin, D. & van der Werf, J.H. 2020. Use of gene expression and whole-  
427 genome sequence information to improve the accuracy of genomic prediction for carcass  
428 traits in Hanwoo cattle. *Genetics Selection Evolution*, 52(1):1-16.
- 429 Erbe, M., Hayes, B., Matukumalli, L., Goswami, S., Bowman, P., Reich, C., Mason, B. &  
430 Goddard, M. 2012. Improving accuracy of genomic predictions within and between dairy  
431 cattle breeds with imputed high-density single nucleotide polymorphism panels. *Journal of  
432 dairy science*, 95(7):4114-4129.

- 433 Fang, L., Cai, W., Liu, S., Canela-Xandri, O., Gao, Y., Jiang, J., Rawlik, K., Li, B.,  
434 Schroeder, S.G. & Rosen, B.D. 2020. Comprehensive analyses of 723 transcriptomes  
435 enhance genetic and biological interpretations for complex traits in cattle. *Genome research*,  
436 30(5):790-801.
- 437 Fang, L., Sahana, G., Ma, P., Su, G., Yu, Y., Zhang, S., Lund, M.S. & Sørensen, P. 2017a.  
438 Exploring the genetic architecture and improving genomic prediction accuracy for mastitis  
439 and milk production traits in dairy cattle by mapping variants to hepatic transcriptomic  
440 regions responsive to intra-mammary infection. *Genetics Selection Evolution*, 49(1):1-18.
- 441 Fang, L., Sahana, G., Ma, P., Su, G., Yu, Y., Zhang, S., Lund, M.S. & Sørensen, P. 2017b.  
442 Use of biological priors enhances understanding of genetic architecture and genomic  
443 prediction of complex traits within and between dairy cattle breeds. *BMC genomics*,  
444 18(1):604.
- 445 Fink, T., Lopdell, T.J., Tiplady, K., Handley, R., Johnson, T.J., Spelman, R.J., Davis, S.R.,  
446 Snell, R.G. & Littlejohn, M.D. 2020. A new mechanism for a familiar mutation–bovine  
447 DGAT1 K232A modulates gene expression through multi-junction exon splice enhancement.  
448 *BMC genomics*, 21(1):1-13.
- 449 Fuchsberger, C., Abecasis, G.R. & Hinds, D.A. 2014. minimac2: faster genotype imputation.  
450 *Bioinformatics*, 31(5):782-784.
- 451 Hayes, B.J. & Daetwyler, H.D. 2018. 1000 Bull Genomes Project to Map Simple and  
452 Complex Genetic Traits in Cattle: Applications and Outcomes. *Annual review of animal*  
453 *biosciences*.
- 454 Howie, B., Fuchsberger, C., Stephens, M., Marchini, J. & Abecasis, G.R. 2012. Fast and  
455 accurate genotype imputation in genome-wide association studies through pre-phasing.  
456 *Nature genetics*, 44(8):955.
- 457 Kern, C., Wang, Y., Xu, X., Pan, Z., Halstead, M., Chanthavixay, G., Saelao, P., Waters, S.,  
458 Xiang, R. & Chamberlain, A. 2021. Functional annotations of three domestic animal  
459 genomes provide vital resources for comparative and agricultural research. *Nature*  
460 *Communications*, 12(1):1-11.
- 461 Koufariotis, L.T., Chen, Y.-P.P., Stothard, P. & Hayes, B.J. 2018. Variance explained by  
462 whole genome sequence variants in coding and regulatory genome annotations for six dairy  
463 traits. *BMC genomics*, 19(1):237.
- 464 Lee, S.H. & Van der Werf, J.H. 2016. MTG2: an efficient algorithm for multivariate linear  
465 mixed model analysis based on genomic information. *Bioinformatics*, 32(9):1420-1422.
- 466 Li, Y.I., van de Geijn, B., Raj, A., Knowles, D.A., Petti, A.A., Golan, D., Gilad, Y. &  
467 Pritchard, J.K. 2016. RNA splicing is a primary link between genetic variation and disease.  
468 *Science*, 352(6285):600-604.
- 469 Liu, S., Fang, L., Zhou, Y., Santos, D.J.A., Xiang, R., Daetwyler, H.D., Chamberlain, A.J.,  
470 Cole, J.B., Li, C.J., Yu, Y., Ma, L., Zhang, S. & Liu, G.E. 2019. Analyses of inter-individual  
471 variations of sperm DNA methylation and their potential implications in cattle. *BMC*  
472 *Genomics*, 20(1):888.
- 473 Loh, P.-R., Danecek, P., Palamara, P.F., Fuchsberger, C., Reshef, Y.A., Finucane, H.K.,  
474 Schoenherr, S., Forer, L., McCarthy, S. & Abecasis, G.R. 2016. Reference-based phasing  
475 using the Haplotype Reference Consortium panel. *Nature genetics*, 48(11):1443.
- 476 Lopdell, T.J., Tiplady, K., Struchalin, M., Johnson, T.J., Keehan, M., Sherlock, R., Couldrey,  
477 C., Davis, S.R., Snell, R.G. & Spelman, R.J. 2017. DNA and RNA-sequence based GWAS  
478 highlights membrane-transport genes as key modulators of milk lactose content. *BMC*  
479 *genomics*, 18(1):1-18.
- 480 MacLeod, I., Bowman, P., Vander Jagt, C., Haile-Mariam, M., Kemper, K., Chamberlain, A.,  
481 Schrooten, C., Hayes, B. & Goddard, M. 2016. Exploiting biological priors and sequence

482 variants enhances QTL discovery and genomic prediction of complex traits. *BMC genomics*,  
483 17(1):144.

484 McLaren, W., Gil, L., Hunt, S.E., Riat, H.S., Ritchie, G.R.S., Thormann, A., Flicek, P. &  
485 Cunningham, F. 2016. The Ensembl Variant Effect Predictor. *Genome Biology*, 17(1):122.

486 Prowse-Wilkins, C., Wang, J., Xiang, R., Goddard, M. & Chamberlain, A. 2021. Putative  
487 causal variants are enriched in annotated functional regions from 6 bovine tissues. *Submitted*.

488 Rosen, B.D., Bickhart, D.M., Schnabel, R.D., Koren, S., Elsik, C.G., Tseng, E., Rowan, T.N.,  
489 Low, W.Y., Zimin, A. & Couldrey, C. 2020. De novo assembly of the cattle reference  
490 genome with single-molecule sequencing. *GigaScience*, 9(3):giaa021.

491 Sargolzaei, M., Chesnais, J.P. & Schenkel, F.S. 2014. A new approach for efficient genotype  
492 imputation using information from relatives. *BMC Genomics*, 15(1):478.

493 Silva, D.B., Fonseca, L.F., Pinheiro, D.G., Magalhães, A.F., Muniz, M.M., Ferro, J.A., Baldi,  
494 F., Chardulo, L.A., Schnabel, R.D. & Taylor, J.F. 2020. Spliced genes in muscle from Nelore  
495 Cattle and their association with carcass and meat quality. *Scientific reports*, 10(1):1-13.

496 Weissbrod, O., Hormozdiari, F., Benner, C., Cui, R., Ulirsch, J., Gazal, S., Schoech, A.P.,  
497 Van De Geijn, B., Reshef, Y. & Márquez-Luna, C. 2020. Functionally informed fine-  
498 mapping and polygenic localization of complex trait heritability. *Nature Genetics*,  
499 52(12):1355-1363.

500 Xiang, R., Berg, I.v.d., MacLeod, I.M., Hayes, B.J., Prowse-Wilkins, C.P., Wang, M.,  
501 Bolormaa, S., Liu, Z., Rochfort, S.J., Reich, C.M., Mason, B.A., Vander Jagt, C.J.,  
502 Daetwyler, H.D., Lund, M.S., Chamberlain, A.J. & Goddard, M.E. 2019. Quantifying the  
503 contribution of sequence variants with regulatory and evolutionary significance to 34 bovine  
504 complex traits. *Proceedings of the National Academy of Sciences*, 116(39):19398-19408.

505 Xiang, R., Hayes, B.J., Vander Jagt, C.J., MacLeod, I.M., Khansefid, M., Bowman, P.J.,  
506 Yuan, Z., Prowse-Wilkins, C.P., Reich, C.M., Mason, B.A., Garner, J.B., Marett, L.C., Chen,  
507 Y., Bolormaa, S., Daetwyler, H.D., Chamberlain, A.J. & Goddard, M.E. 2018. Genome  
508 variants associated with RNA splicing variations in bovine are extensively shared between  
509 tissues. *BMC Genomics*, 19(1):521.

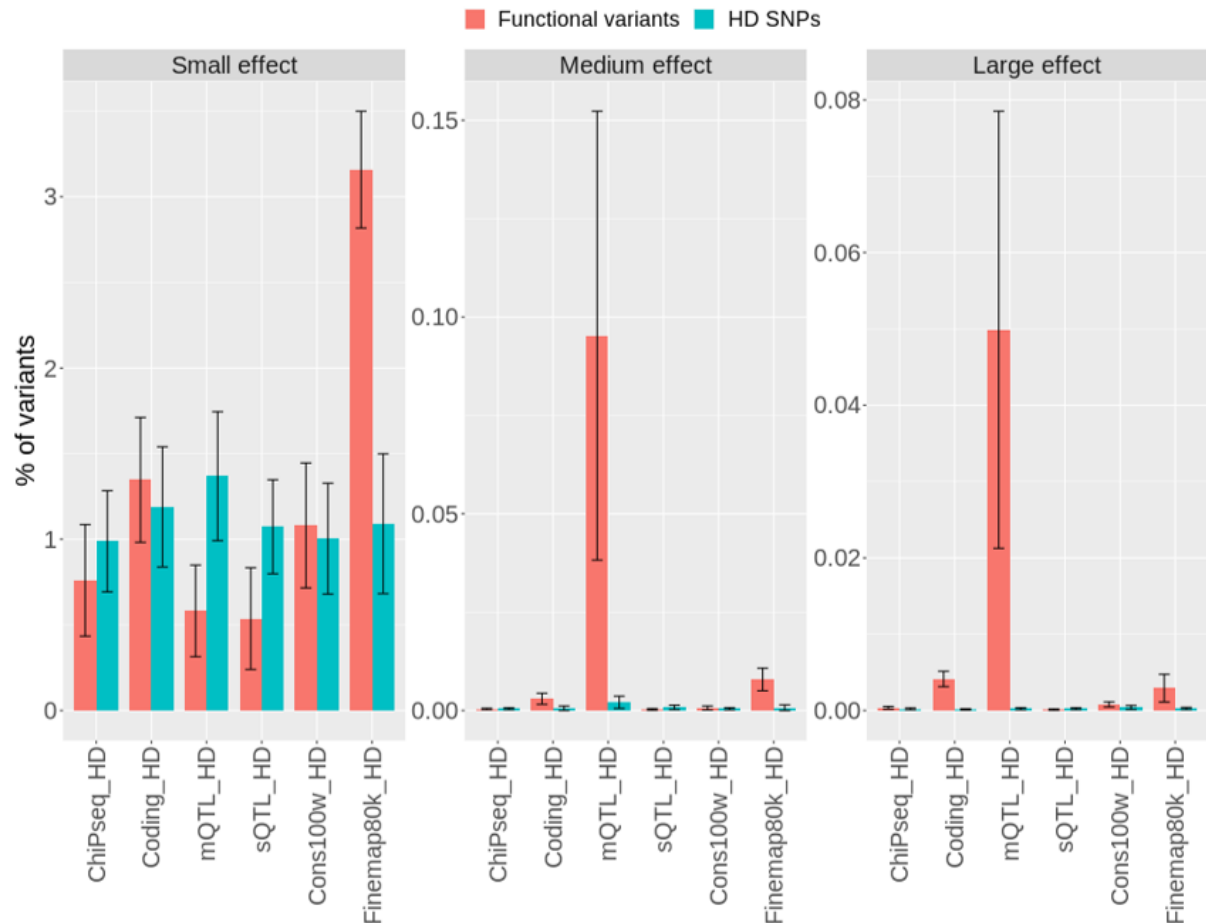
510 Xiang, R., MacLeod, I.M., Daetwyler, H.D., de Jong, G., O'Connor, E., Schrooten, C.,  
511 Chamberlain, A.J. & Goddard, M.E. 2021. Genome-wide fine-mapping identifies pleiotropic  
512 and functional variants that predict many traits across global cattle populations. *Nature*  
513 *Communications*, 12(1):860.

514 Xiang, R., van den Berg, I., MacLeod, I.M., Daetwyler, H.D. & Goddard, M.E. 2020. Effect  
515 direction meta-analysis of GWAS identifies extreme, prevalent and shared pleiotropy in a  
516 large mammal. *Commun Biol*, 3(1):88.

517 Xu, L., Gao, N., Wang, Z., Xu, L., Liu, Y., Chen, Y., Xu, L., Gao, X., Zhang, L. & Gao, H.  
518 2020. Incorporating Genome Annotation Into Genomic Prediction for Carcass Traits in  
519 Chinese Simmental Beef Cattle. *Frontiers in Genetics*, 11.

520

521



522

523

524

525

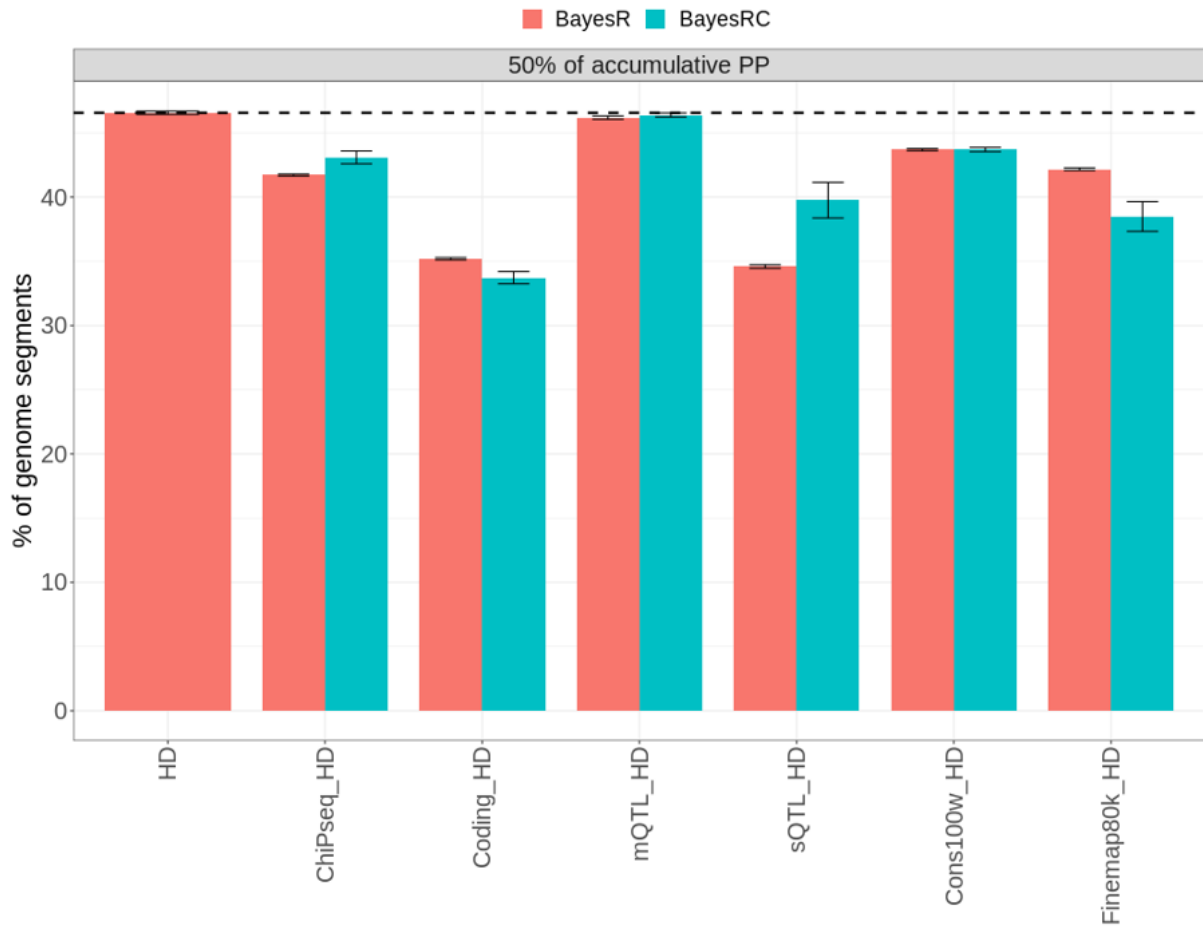
526

527

528

529

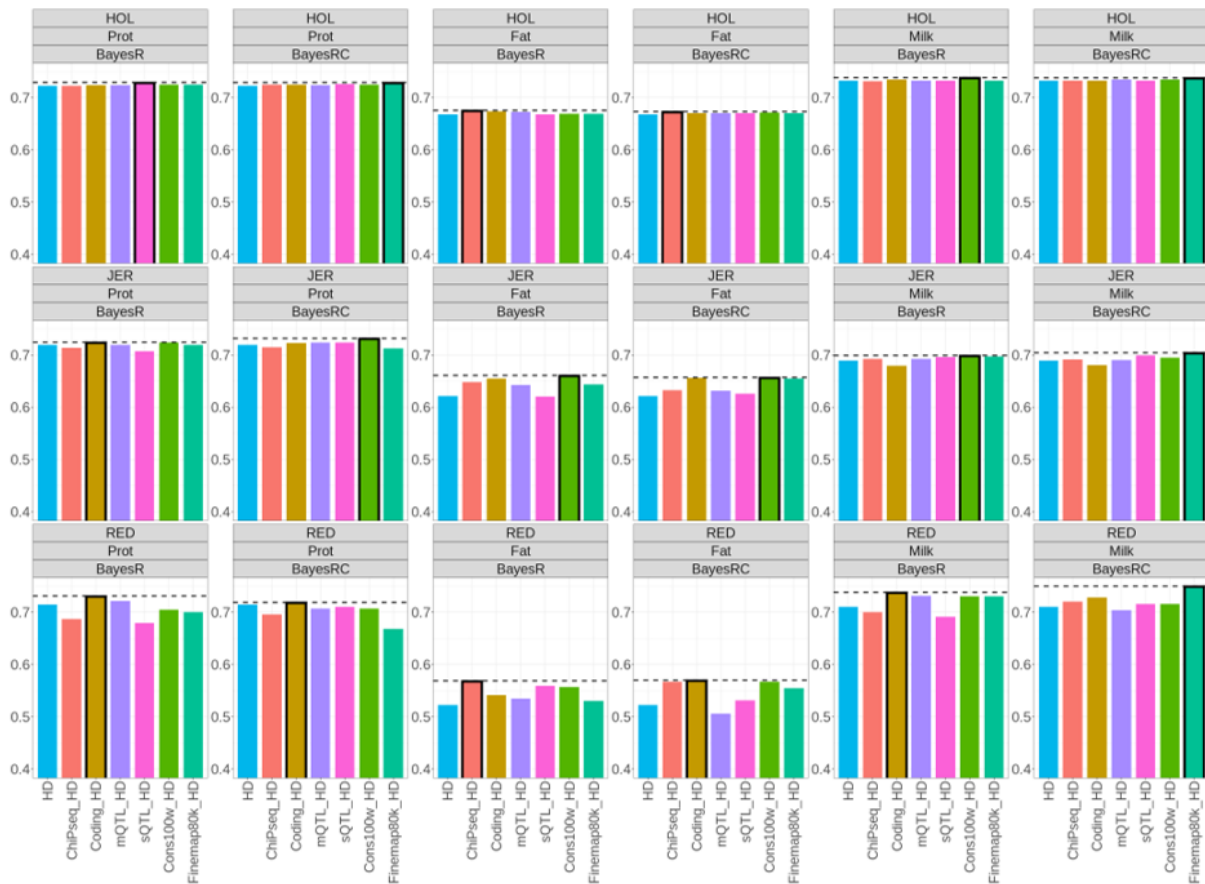
**Figure 1.** The proportion of small-effect, medium effect and large effect variants in functional variants and HD SNPs. The mean and standard error bars are averaged across 6 traits. ChiPseq\_HD: ChIP-seq peaks + HD SNPs. Coding\_HD: coding variants + HD SNPs. mQTL\_HD: mQTLs + HD SNPs. sQTL\_HD: sQTL variants + HD SNPs. Cons100w\_HD: conserved variants across 100 vertebrates + HD SNPs. Finemap80k\_HD: finely mapped variants + HD SNPs.



530

531 **Figure 2.** Mapping precision of different models. The Y-axes represent the percentage of  
532 50kb segments needed by the model to explain 50% of the cumulative sum of posterior  
533 probability (PP) of variants. A shorter bar means less amount of segments the model needs to  
534 explain the same amount of genetic variance, indicating higher mapping precision. Black  
535 dashed line indicates the Y value for the HD SNPs, fitted along in BayesR. ChiPseq\_HD:  
536 ChIP-seq peaks + HD SNPs. Coding\_HD: coding variants + HD SNPs. mQTL\_HD: mQTLs  
537 + HD SNPs. sQTL\_HD: sQTL variants + HD SNPs. Cons100w\_HD: conserved variants  
538 across 100 vertebrates + HD SNPs. Finemap80k\_HD: finely mapped variants + HD SNPs.

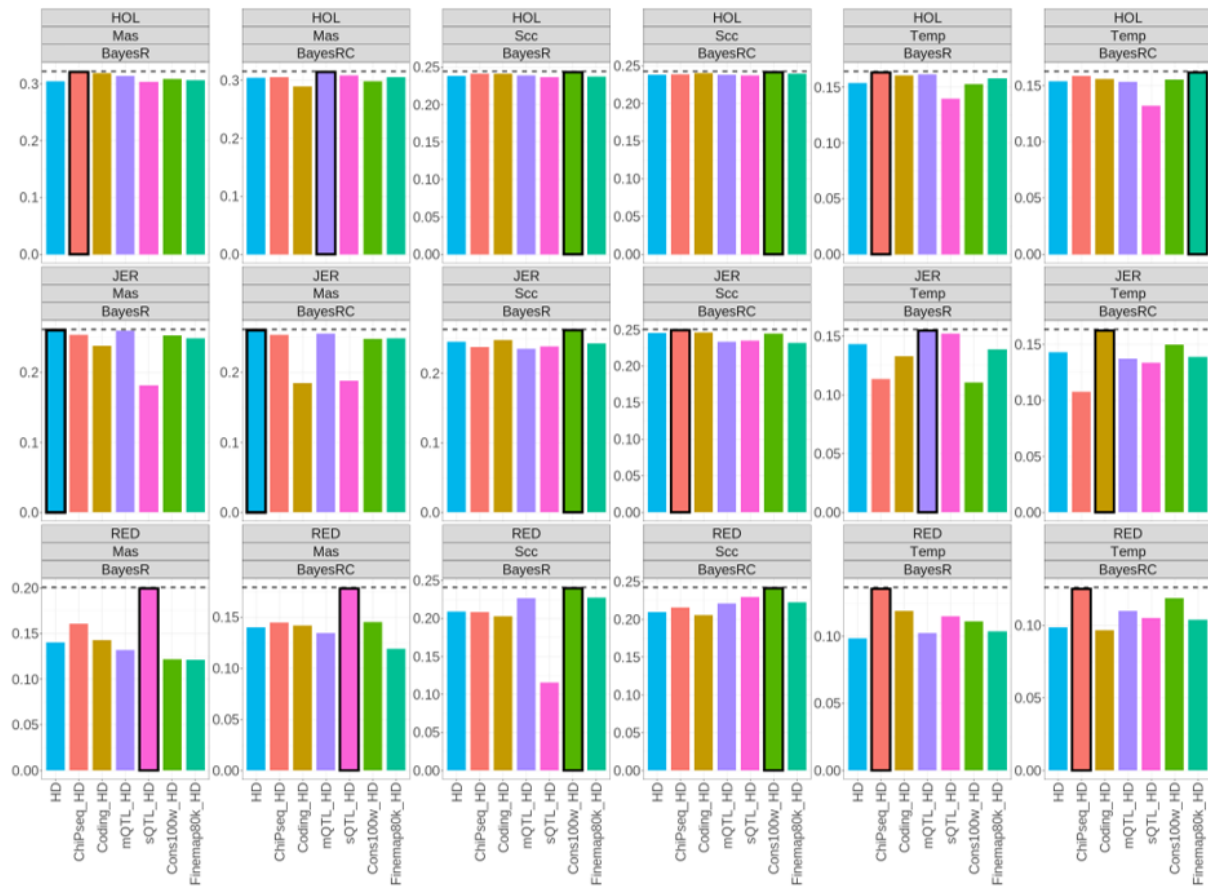
539



540

541 **Figure 3.** Genomic prediction accuracy (Pearson correlation coefficient, Y-axis) for  
 542 production traits, across different functional/evolutionary variant sets, breeds and Bayesian  
 543 methods. A black border and a dashed line of a bar indicate that it has the highest genomic  
 544 prediction accuracy in the panel. HOL: Holstein breed. JER: Jersey breed. RED: Prot: milk  
 545 protein yield. Fat: milk fat yield. Milk: milk yield. Australian Red. ChiPseq\_HD: ChIP-seq  
 546 peaks + HD SNPs. Coding\_HD: coding variants + HD SNPs. mQTL\_HD: mQTLs + HD  
 547 SNPs. sQTL\_HD: sQTL variants + HD SNPs. Cons100w\_HD: conserved variants across 100  
 548 vertebrates + HD SNPs. Finemap80k\_HD: finely mapped variants + HD SNPs.

549



550

551 **Figure 4.** Genomic prediction accuracy (Pearson correlation coefficient, Y-axis) for mastitis,  
 552 somatic cell count and temperament across different functional/evolutionary variant sets,  
 553 breeds and Bayesian methods. A black border and a dashed line of a bar indicate that it has  
 554 the highest genomic prediction accuracy in the panel. HOL: Holstein breed. JER: Jersey  
 555 breed. RED: Australian Red. Mas: mastitis. Scc: somatic cell count. Temp: temperament.  
 556 ChiPseq\_HD: ChIP-seq peaks + HD SNPs. Coding\_HD: coding variants + HD SNPs.  
 557 mQTL\_HD: mQTLs + HD SNPs. sQTL\_HD: sQTL variants + HD SNPs. Cons100w\_HD:  
 558 conserved variants across 100 vertebrates + HD SNPs. Finemap80k\_HD:  
 559 variants + HD SNPs.