

Bayesian genome-wide analysis of cattle traits using variants with functional and evolutionary significance

Ruidong Xiang^{id A,B,C}, Ed J. Breen^B, Claire P. Prowse-Wilkins^{A,B}, Amanda J. Chamberlain^B and Michael E. Goddard^{A,B}

^AFaculty of Veterinary and Agricultural Science, The University of Melbourne, 142 Royal Parade, Parkville, Vic. 3052, Australia.

^BAgriculture Victoria, AgriBio, Centre for AgriBiosciences, 5 Ring Road, Bundoora, Vic. 3083, Australia.

^CCorresponding author. Email: ruidong.xiang@unimelb.edu.au

Abstract

Context. Functional genomics studies have highlighted genomic regions with regulatory and evolutionary significance. Such information independent of association analysis may benefit fine-mapping and genomic selection of economically important traits. However, systematic evaluation of the use of functional information in mapping, and genomic selection of cattle traits, is lacking. Also, single-nucleotide polymorphisms (SNPs) from the high-density (HD) panel are known to tag informative variants, but the performance of genomic prediction using HD SNPs together with variants supported by different functional genomics is unknown.

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

Methods. Two models were used, namely (1) BayesR, which includes priors of four distribution of variant effects, and (2) BayesRC, which includes additional priors of different functional classes of variants. Bayesian models were trained in three breeds of 28 000 cows of Holstein, Jersey and Australian Red and predicted into 2600 independent bulls.

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

Conclusions. Our findings demonstrated the usefulness of functional data in genomic mapping and prediction.

Implications. We have highlighted the need for effective tools exploiting complex functional datasets to improve genomic prediction.

Additional keywords: functional genomics, animal breeding, genetic mapping, quantitative genetics.

Received 8 February 2021, accepted 12 May 2021, published online 21 July 2021

Introduction

Emerging evidence shows that genomic variants with causal roles in biology can be used to improve genomic prediction of complex traits. The biological function of genomic variants provides information independent of genotype-trait associations that are usually confounded by linkage disequilibrium (LD). Such independent information can be exploited to identify informative variants. Once identified, informative variants can be used to improve genomic prediction (Xiang *et al.* 2021). While the use of functional data in improving genomic mapping and prediction has

been reported in humans (Amariuta *et al.* 2020; Weissbrod *et al.* 2020), using functional data in predicting the genetic merit of animal traits has not been comprehensively examined. However, there is evidence in cattle supporting the advantage of the use of functional information in genomic mapping and prediction with the linear mixed model (Fang *et al.* 2017a, 2017b; Liu *et al.* 2019; Xiang *et al.* 2019; Xu *et al.* 2020).

The Functional Annotation of ANimal Genomes (FAANG) consortium (Clark *et al.* 2020) provides many types of sequencing data indicating the functionality of

genome-wide sites (examples reviewed in Clark *et al.* 2020). While these public datasets await exploitation, the structure and information content of different functional datasets vary significantly. For example, we recently showed that among all analysed functional datasets, a set of 300 000+ sequence variants within sites highly conserved across 100 vertebrate species had the strongest enrichment with cattle trait heritability (Xiang *et al.* 2019), which primarily influences genomic prediction accuracy. Additionally, a few thousand variants affecting the concentration of milk fat metabolites, i.e. metabolic quantitative trait loci (mQTLs), also had significantly higher variance than did single-nucleotide polymorphisms (SNPs) in the 50 K panel for cattle traits. Millions of variants that change gene-expression levels (geQTLs) or RNA splicing (sQTLs) are also enriched with complex-trait QTL (Li *et al.* 2016; Lopdell *et al.* 2017; Xiang *et al.* 2018; Fink *et al.* 2020; Silva *et al.* 2020). However, recent studies have shown that variants close to genes with high or specific expression patterns have limited improvement in prediction accuracy (de las Heras-Saldana *et al.* 2020; Fang *et al.* 2020). Another common type of functional data are peaks from ChIP-seq for histone modifications, which are enriched with promoters and/or enhancers regulating gene activities (Carey *et al.* 2009). Our work showed that hundreds of thousands of variants under ChIP-seq peaks are enriched for complex-trait QTL in cattle (Xiang *et al.* 2019; Prowse-Wilkins *et al.* 2021). In addition, variants within the gene-coding regions are expected to have a high impact on complex traits. However, we and others previously found that coding-related variants (~100 000) have limited contributions to cattle trait heritability (Koufariotis *et al.* 2018; Xiang *et al.* 2019), although their use in improving genomic prediction has not been studied.

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

In the present study, we evaluate sequence-variant sets prioritised by six types of functional and evolutionary data in combination with the standard HD SNPs in genomic mapping and prediction of six dairy cattle traits. We train the prediction equations by using the BayesR method (Erbe *et al.* 2012), which fits a mixture of four distributions of variant effects, and by using the BayesRC method, which fits different

distributions for each functional class of variant classifications (MacLeod *et al.* 2016). Genomic predictors were trained using 28 000 cows that included three breeds, namely, Holstein, Jersey and Australian Red. Genomic estimated breeding values (gEBVs) were predicted and validated in 2500 Holstein, Jersey and Australian Red bulls. We compare the results of mapping and genomic prediction across the above-described scenarios, discuss these results and provide suggestions for future studies.

Materials and methods

The phenotype data analysed in the present study were collected by DataGene Australia (<http://www.datagene.com.au/>) and no further live-animal experimentation was required for our analyses. A set of 28 049 Australian cows was used as the discovery population and a set of 2567 bulls was used as the validation population. The bull phenotypes were obtained as daughter trait deviations, that is, the average trait deviations of a bull's daughters pre-corrected for known fixed effects by DataGene. The cow phenotypes were measured on themselves. Note that these bulls and cows were not included in the 44 000+ animals used to discover functional variants (Xiang *et al.* 2019, 2020, 2021). We also checked the pedigree to make sure that bulls used in the validation population were not the sires of cows from the discovery population. Cows in the discovery set included 24 305 Holstein, 2486 Jersey, and 1258 Australian Red. Bulls in the validation datasets contained 2091 Holstein, 385 Jersey and 91 Australian Red. Traits considered in the analysis included protein yield, fat yield, milk yield, mastitis (Mas), somatic cell count (Scc) and temperament (Temp).

The genotypes used in the study were imputed sequence variants based on Run7 of the 1000-Bull Genomes Project (Daetwyler *et al.* 2014; Hayes and Daetwyler 2018), based on the ARS-UCD1.2 reference bovine genome (https://www.ncbi.nlm.nih.gov/assembly/GCF_002263795.1/; Rosen *et al.* 2020). Variants with Minimac3 (Howie *et al.* 2012; Fuchsberger *et al.* 2015) imputation accuracy $R^2 > 0.4$ and minor allele frequency of >0.005 in bulls and cows. Most bulls were genotyped with a medium-density SNP array (50 K) or a high-density SNP array and most cows were genotyped with a low-density panel of ~6900 SNPs overlapping the standard-50 K panel (BovineSNP50 beadchip, Illumina Inc.). The low-density genotypes were first imputed to the Standard-50 K panel and then all 50 K genotypes were imputed to the HD panel using Fimpute v3 (Sargolzaei *et al.* 2014; Xiang *et al.* 2019). Then, all HD genotypes were imputed to sequence using Minimac3 with Eagle (v2) to pre-phase genotypes (Howie *et al.* 2012; Loh *et al.* 2016).

We aimed to test whether variant sets selected from different functional and/or evolutionary information, in addition to the standard HD SNP panel, can be useful for genomic prediction. Therefore, we first included a baseline set, which is 610 764 SNPs from the standard bovine HD panel. There were six functional and/or evolutionary variant sets: 549 007 variants under multiple ChIP-seq peaks (Kern *et al.* 2021; Prowse-Wilkins *et al.* 2021; ChIPseq), 106 538 variants annotated as related to coding activities by Ensembl Variant

Effect Predictor (McLaren *et al.* 2016; Coding), 943 315 variants affecting RNA splicing sQTLs from four cattle tissues (Chamberlain *et al.* 2018; Xiang *et al.* 2018; Daetwyler *et al.* 2019; sQTL), 65 394 finely mapped variants with pleiotropic effects genome-wide (Xiang *et al.* 2021; Finemap80k), 4871 variants affecting milk fat metabolite mQTLs (Xiang *et al.* 2019; mQTL) and 317 279 conserved sites across 100 vertebrates (Xiang *et al.* 2019; Cons100w). Note that some of these functional variant sets were initially determined on the UMD3.1 genome and were from different cattle populations. These sets were lifted over from the older genome to ARS-UCD1.2 and filtered with imputation accuracy and minor allele frequency in the new cattle populations.

The model training of the above-described data used BayesR (Erbe *et al.* 2012) and BayesRC (MacLeod *et al.* 2016), which are now implemented via BayesR3, with improved efficiency using blocks. BayesR jointly models all variants together, with different effect distribution priors. BayesRC follows the same approach but, in addition, allows a C prior which models classes of variants. Another aim is to see whether there are differences in genomic prediction accuracy by modelling the same variants by using BayesR and BayesRC. To aid this comparison, we combined each functional variant set with the HD variants, which led to the following six combined variant sets: (1) ChIP-seq peak-tagged variants + HD SNPs (ChIPseq_HD), (2) coding variants + HD variants (Coding_HD), (3) sQTL variants + HD SNPs (sQTL_HD), (4) finely mapped variants + HD SNPs ('Finemap80k_HD'), (5) mQTL variants + HD SNPs (mQTL_HD) and (6) conserved variants + HD SNPs (Cons100w_HD). The average minor allele frequency of these sets of variants was 0.22 (± 0.00014) for ChIPseq_HD, 0.25 (± 0.0002) for Coding_HD, 0.24 (± 0.0001) for sQTL_HD, 0.27 (± 0.0002) for Finemap80k_HD, 0.27 (± 0.0002) for mQTL_HD, 0.23 (± 0.0002) for Cons100w_HD, and 0.27 (± 0.0002) for HD alone.

In single-trait BayesR, we directly model these six variant sets one set at a time. To create a reference baseline, we also used single-trait BayesR to fit the HD variant set (HD) alone. In single-trait BayesRC, for each of the same six combined variant sets, we specified the following two different variant classes: (1) variants appeared in the functional and/or evolutionary set and (2) variants appeared only in the HD variant set.

Both BayesR and BayesRC modelled variant effects as a mixture distribution of four normal distributions, including a null distribution, $N(0, 0.0\sigma_g^2)$, and three others, namely, $N(0, 0.0001\sigma_g^2)$, $N(0, 0.001\sigma_g^2)$ and $N(0, 0.01\sigma_g^2)$, where σ_g^2 is the additive genetic variance for the trait. The starting value of σ_g^2 for each trait was estimated using GREML implemented in the MTG2 (Lee and Van der Werf 2016), with a single genomic relationship matrix made of all sequence variants. The statistical model used in the single-trait BayesR and BayesRC was

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

where \mathbf{y} is a vector of phenotypic records, \mathbf{W} is the design matrix of marker genotypes, centred and standardised to have a

unit variance, \mathbf{v} is the vector of variant effects, distributed as a mixture of the four distributions as described above; \mathbf{X} is the design matrix allocating phenotypes to fixed effects; \mathbf{b} is the vector of fixed effects, including breeds; \mathbf{e} is vector of residual errors. As a result, the effect b for each variant jointly estimated with other variants was obtained for further analysis.

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

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

$$gEBV = Z\hat{s} \quad (2)$$

where Z is a matrix of the standardised genotypes of animals in the validation set, and \hat{s} is the vector of variant effects from the training model. In addition, to test whether adding a large number of random variants to the HD panel can increase mapping precision and prediction accuracy, a random set of 944 616 variants matching the size of the largest set of functional variants (sQTL, 943 315 variants) was also selected and added to the HD panel (Random_HD). This random set was analysed for BayesR, mapping precision and prediction accuracy in the same fashion as were other variant sets described above.

Results

Information content in the functional variant sets

Averaged across mixing proportions from single-trait BayesRC, we show that compared with HD SNPs, the finely mapped variants had consistently higher enrichment, with variants showing small, medium and large effects (Fig. 1). Variants within coding regions showed higher enrichment than did HD SNPs for large- and medium-effect variants. Interestingly, mQTLs, which were variants affecting the concentration of milk-fat metabolites (Benedet *et al.* 2019; Xiang *et al.* 2019), had lower enrichment of small-effect

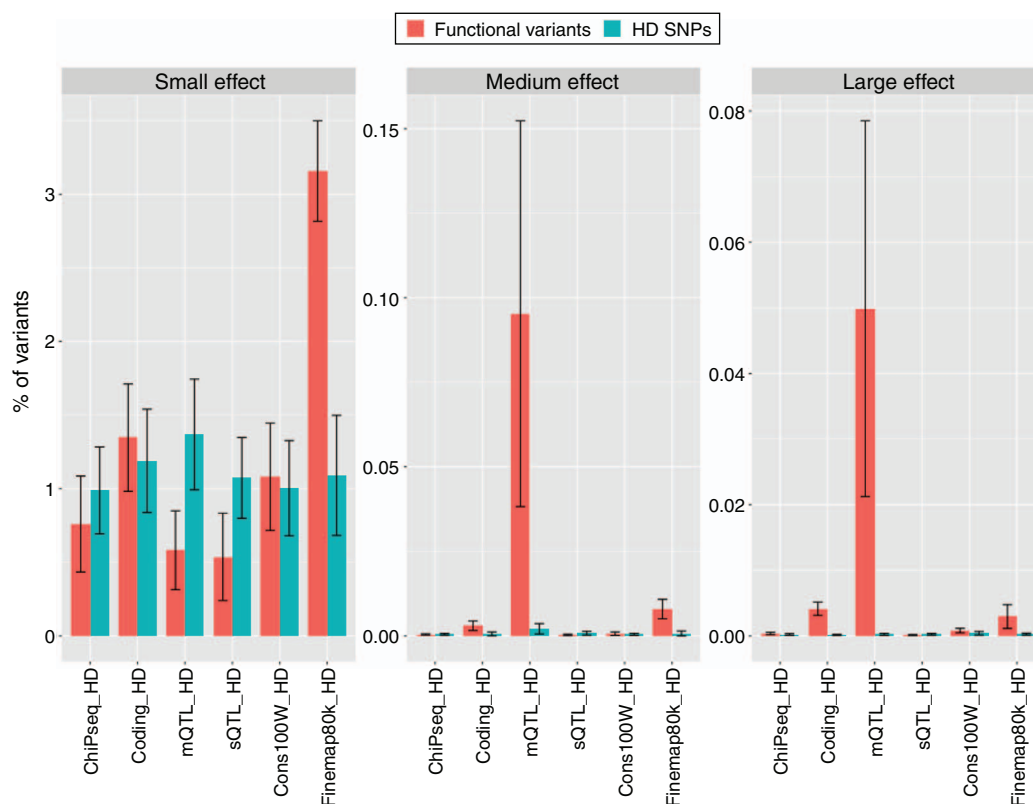


Fig. 1. The proportion of small-effect, medium-effect and large-effect variants in functional variants and HD SNPs. The means and standard error bars are averaged across six traits. ChiPseq_HD, ChIP-seq peaks + HD SNPs. Coding_HD, coding variants + HD SNPs. mQTL_HD, mQTLs + HD SNPs. sQTL_HD, sQTL variants + HD SNPs. Cons100w_HD, conserved variants across 100 vertebrates + HD SNPs. Finemap80k_HD, finely mapped variants + HD SNPs.

variants than did HD SNPs, but had higher enrichment of medium- and large-effect variants than did HD SNPs.

Mapping precision

Across traits, we have shown that all models using functional variants, except mQTL, needed a smaller amount of genome-wide segments to explain 50% of the cumulative sum of PP, than did HD SNPs (Fig. 2). This means that when adding to the HD SNPs, most functional variants increased mapping precision. In contrast, adding randomly selected 944 000 variants to HD SNPs increased the amount of genome-wide segments (by $2.82\% \pm 0.13\%$) across scenarios to explain 50% of the cumulative sum of PP, compared with using only HD SNPs. This suggested that adding random variants to HD decreases mapping precision. It is worth noting that when using 106 538 coding variants and 65 394 finely mapped variants, BayesRC provided a further increase in mapping precision over HD SNPs from BayesR. In contrast, when using 549 007 ChIP-seq-tagged variants and 943 315 sQTL variants, BayesRC had less increase in mapping precision over HD SNPs than did BayesR. This could be due to the reduced signal-to-noise ratio in large variant sets of ChIP-seq-tagged variants and sQTLs.

Genomic prediction of traits

In total, we evaluated the genomic prediction accuracy in 216 scenarios, across six single-trait analysis, six functional categories, four breeds in the validation population, and two Bayesian methods. Out of these 216 scenarios, 142 (66%) times, HD SNPs combined with functional variants increased genomic prediction accuracy, compared with the prediction using only the HD SNPs (Figs 3, 4). In 51 of 216 times (24%), the increase in prediction accuracy ($[r_{\text{functional}} - r_{\text{HD}}] \times 100\%$) was greater than 1%. These 51 cases were almost all accounted for by Jersey (15/51) and Australian Red (34/51), with only two cases in Holstein cattle. In 29 analyses (14%), the increase in prediction accuracy over HD SNPs was greater than 2%. All these 29 cases were for non-Holstein breeds. Among tested functional sets, genomic prediction accuracy was the best when the HD variants were combined with conserved variants (Cons100w_HD). In contrast, averaged across tested scenarios, adding randomly selected 944 000 variants to HD had a slightly worse or no improvement in prediction accuracy ($-0.5\% \pm 0.49\%$) compared with using only the HD panel to predict.

As shown in Fig. 3, the genomic prediction accuracy of milk production traits using HD SNPs in Holstein cattle was

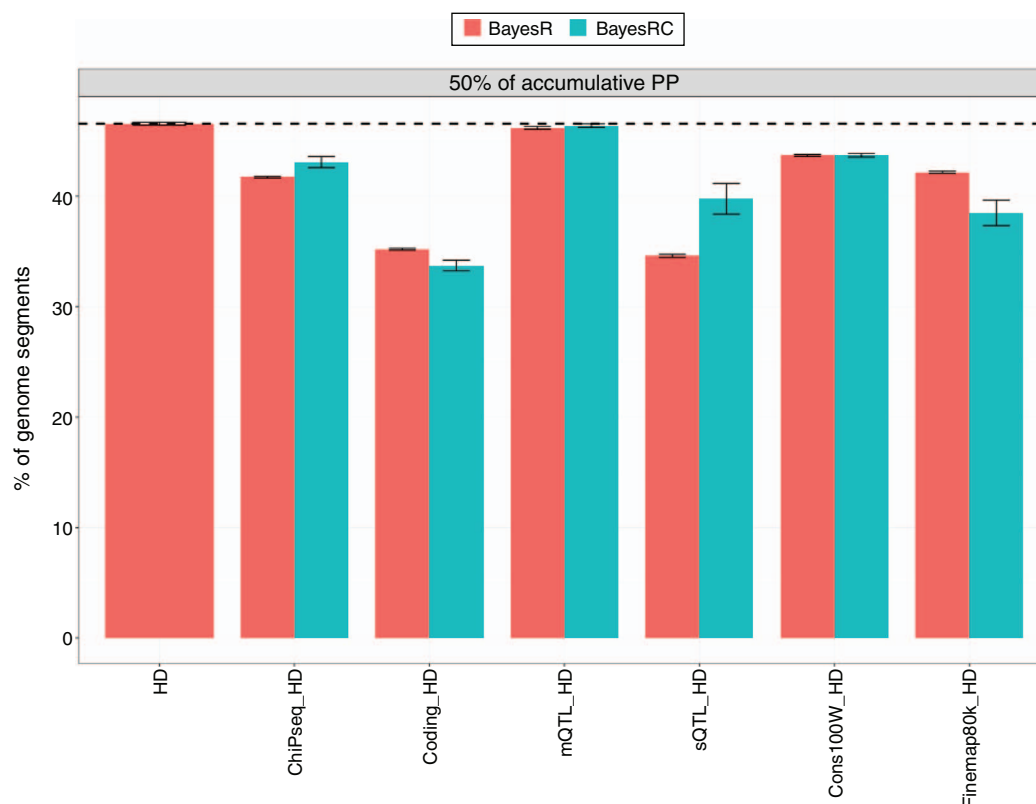


Fig. 2. Mapping precision of different models. The y-axes represent the percentage of 50 kb segments needed by the model to explain 50% of the cumulative sum of posterior probability (PP) of variants. A shorter bar means that the model needs a smaller amount of segment to explain the same amount of genetic variance, indicating higher mapping precision. Black dashed line indicates the Y-value for the HD SNPs, fitted along in BayesR. 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.

already high (~0.7) and the increases in accuracy from functional variants were very small. However, larger increases were evident in Jersey and Australian Red. For milk production traits, 10 times of 18, the genomic prediction accuracy was the most improved by conserved variants and coding variants combined with HD SNPs, followed by finely mapped variants combined with HD SNPs (4/18), ChIP-seq-tagged variants (3/18) combined with HD SNPs. sQTL combined with HD variants had the highest accuracy when predicting protein yield in Holstein.

As shown in Fig. 4, the greatest increases in prediction accuracy for traits Mas, Scc and Temp were again seen in non-Holstein breeds. Chip-seq peak-tagged variants combined with HD SNPs (5/18 times) and conserved variants combined with HD SNPs (5/18 times) had the best performances in predicting Mas, Scc and Temp.

Across all scenarios, we did not see a clear distinction in prediction accuracy between BayesR and BayesRC in the current study. There may be some tendencies where BayesRC had a higher accuracy than did BayesR for Scc, Mas and Temp. However, none of these differences was significant.

Discussion

Our systematic evaluations showed that functional information can improve genomic mapping and prediction of cattle traits, even when HD SNPs are used, although there were times where HD SNPs alone still had robust performances. It is usually the less represented breeds, such as Jersey and Australian Red, that benefit the most from the improvements using functional data. This suggests that functional information can well complement HD SNPs, especially in breeds with smaller training sets. Adding randomly selected variants to the HD panel reduced mapping precision and provided no improvement in prediction accuracy, compared with using only the HD panel. This supports that the benefit provided via selecting variants based on functional importance cannot be achieved simply by adding more sequence variants.

We showed that the biological information content, which can be used to benefit mapping and/or prediction, is different among functional datasets. One of the top-performing functional variant sets in mapping large-effect variants was the finely mapped 80 000 variants (Xiang *et al.* 2021). This

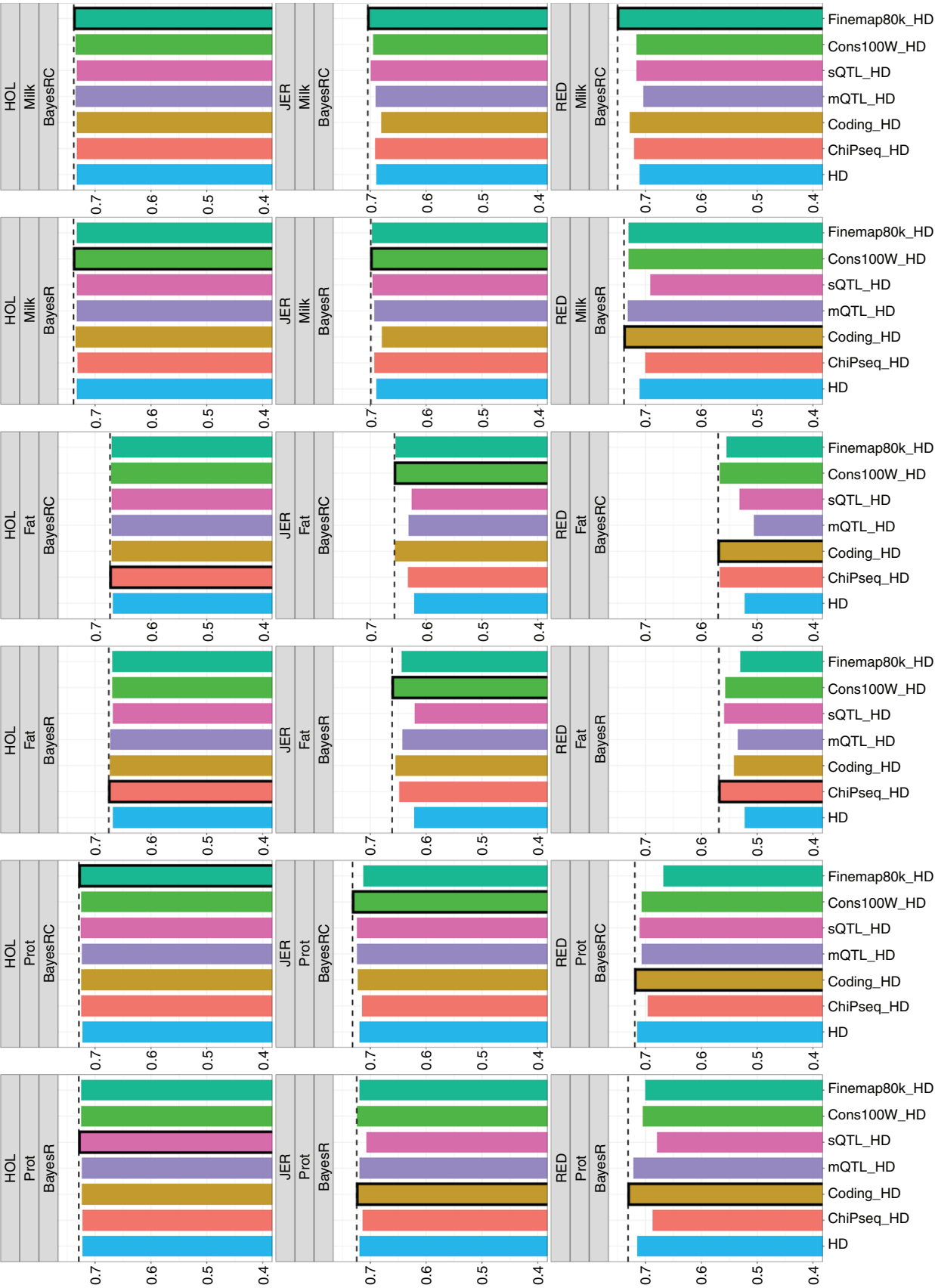


Fig. 3. Genomic prediction accuracy (Pearson correlation coefficient, y-axis) for production traits, across different functional/evolutionary variant sets, breeds and Bayesian methods. A black border and a dashed line of a bar indicate that it has the highest genomic prediction accuracy in the panel. HOL, Holstein breed. JER, Jersey breed. RED, Australian Red. Prot, milk protein yield. Fat, milk fat yield. Milk, milk yield. 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.

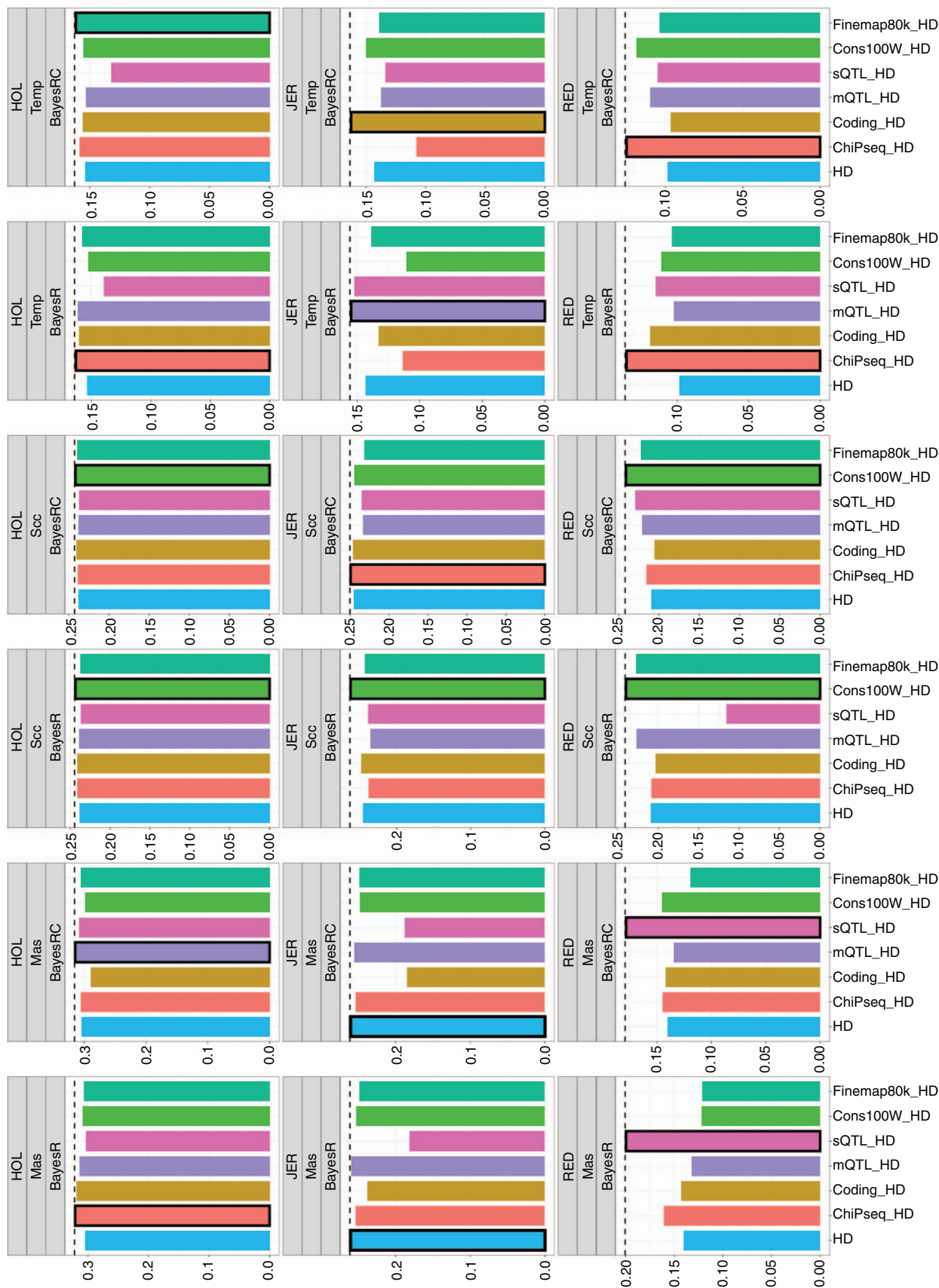


Fig. 4. Genomic prediction accuracy (Pearson correlation coefficient, y-axis) for mastitis, somatic cell count and temperament across different functional/evolutionary variant sets, breeds and Bayesian methods. A black border and a dashed line of a bar indicate that it has the highest genomic prediction accuracy in the panel. HOL, Holstein breed. JER, Jersey breed. RED, Austrian Red. Mas, mastitis. Sec, somatic cell count. Temp, temperament. 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.

result was somewhat expected as these variants combined information from multiple functional datasets and also included variants affecting multiple dairy cattle traits. These finely mapped 80 000 variants outperformed the SNPs from the 50 K panel in previous evaluations (Xiang *et al.* 2021). Furthermore, finely mapped 80 000 variants showed enhanced enrichment of large-effect variants and improvement in mapping precision when modelled with BayesRC. This suggests that this much more refined set of variants (chosen because they were more relevant to the traits of interest) is likely to be more enriched for variants that are more strongly associated with the trait or are causal. BayesRC would only outperform BayesR when there is strong enrichment for QTL in at least one of the defined classes. The other functional groups tested are not trait specific (except mQTL for fat), and so are likely to be less enriched than each trait.

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

Unlike the results in mapping large-effect variants, for genomic prediction, the top-performing variant set is the conserved variants combined with HD SNPs. The advantage of adding conserved variants to HD SNPs was particularly evident when predicting Scc, Mas and Temp of non-Holstein breeds (Fig. 4). In fact, in these scenarios, HD SNPs alone did not perform so well and this leaves more room for functional variants to improve the prediction accuracy. Another variant set that performed well in genomic prediction is the set of ChIP-seq peak-tagged variants. Again, such an advantage was the most evident when predicting Scc, Mas and Temp in non-Holstein breeds. Interestingly, ChIP-seq variants combined with HD SNPs appear to show some particular advantages in predicting Temp. There may be some large-effect variants for Temp captured by ChIP-seq peaks.

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

As mentioned earlier, BayesRC would outperform BayesR only when there is strong enrichment for QTL in at least one of the defined classes. It would also require functional information to be trait-specific. We saw advantages in BayesRC over BayesR in detecting enrichment with large-effect variants by using finely mapped variants, coding variants and mQTLs. BayesRC also had advantages over BayesR in mapping precision when used with finely mapped variants and coding variants. While these functional data are expected to be informative, they did not provide consistent advantages for BayesRC to predict traits over BayesR. Across all tested cases, we did not see strong advantages in BayesRC over BayesR in genomic prediction (Fig. 4). BayesRC may have some tendencies to better predict Scc, Mas and Temp than does BayesR. However, the differences were not statistically significant. The reason behind these observations may be complex.

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

We therefore suggest two future research directions to improve on the current results. The first is to increase the information content in functional datasets. This can be achieved by either increasing the sample size (biological replicates, tissues and experimental conditions) of functional datasets or by developing better bioinformatic tools to increase the signal-to-noise ratio in functional datasets before they can be processed by genomic prediction models. The second direction is to improve the current genomic prediction models. Because the type and complexity of functional data will keep growing, it will be necessary to develop more sophisticated and flexible methods to better extract information from complex functional data. For example, an extended BayesRC that can model quantitative biological priors, instead of qualitative classes, will be needed. Similarly, in the future we will use larger sample sizes and diverse breeds in the training model to reduce LD between sequence variants. This will also increase the need for Bayesian methods to be more efficient.

In conclusion, our evaluation of Bayesian genomic prediction using functional and evolutionary information with HD SNPs provides novel insights into this emerging area. We show that functional datasets of conserved variants, coding variants, ChIP-seq peaks and previously finely mapped variants can improve genomic mapping and/or genomic prediction, even when HD SNPs are used. Such improvements usually benefit non-Holstein breeds, given the current available functional datasets. We found that by

using informative biological priors, BayesRC has significant advantages over BayesR in detecting enrichment with large-effect variants and in mapping precision. However, the advantage of BayesRC over BayesR for genomic prediction was not consistent. Our results highlighted the need to develop better tools to extract information from complex functional datasets, which will benefit genomic prediction in large datasets. Fusing functional genomics with genomic selection presents great opportunities to develop new technologies that improve animal breeding and genetics.

Conflicts of interest

The authors declare no conflicts of interest.

Declaration of funding

Australian Research Council's Discovery Projects (DP160101056 and DP200100499) supported R. X. and M. E. G. DairyBio, a joint venture project between Agriculture Victoria (Melbourne, Australia), Dairy Australia (Melbourne, Australia) and the Gardiner Foundation (Melbourne, Australia), funded computing resources used in the analysis. The authors also thank the University of Melbourne, Australia, for supporting this research. No funding bodies participated in the design of the study nor analysis, or interpretation of data nor in writing the manuscript.

Acknowledgements

We thank DataGene and CRV who provided access to the reference data used in this study. We thank Gert Nieuwhof, Kon Konstantinov and Timothy P. Hancock (DataGene) and staff from DairyNZ for the preparation and provision of data. We thank Dr Sunduimijid Bolormaa for the sequence variant data imputation. We thank Drs Iona M. MacLeod and Hans D. Daetwyler for critical reading of the manuscript.

References

- Amariuta T, Ishigaki K, Sugishita H, Ohta T, Koido M, Dey KK, Matsuda K, Murakami Y, Price AL, Kawakami E (2020) Improving the trans-ancestry portability of polygenic risk scores by prioritizing variants in predicted cell-type-specific regulatory elements. *Nature Genetics* **52**(12), 1346–1354. doi:10.1038/s41588-020-00740-8
- Benedet A, Ho P, Xiang R, Bolormaa S, De Marchi M, Goddard M, Pryce J (2019) The use of mid-infrared spectra to map genes affecting milk composition. *Journal of Dairy Science* **102**(8), 7189–7203. doi:10.3168/jds.2018-15890
- Carey MF, Peterson CL, Smale ST (2009) Chromatin immunoprecipitation (ChIP). *Cold Spring Harbor Protocols* **4**(9), pdb.prot5279.
- Chamberlain A, Hayes B, Xiang R, Vander Jagt C, Reich C, Macleod I, Prowse-Wilkins C, Mason B, Daetwyler H, Goddard M (2018) Identification of regulatory variation in dairy cattle with RNA sequence data. In '11th World Congress on Genetics Applied to Livestock Production (WCGALP)', Auckland, New Zealand. p. 254.
- Clark EL, Archibald AL, Daetwyler HD, Groenen MA, Harrison PW, Houston RD, Kühn C, Lien S, Macqueen DJ, Reecy JM (2020) From FAANG to fork: application of highly annotated genomes to improve farmed animal production. *Genome Biology* **21**(1), 285. doi:10.1186/s13059-020-02197-8
- Daetwyler HD, Capitan A, Pausch H, Stothard P, Van Binsbergen R, Brøndum RF, Liao X, Djari A, Rodriguez SC, Grohs C (2014) Whole-genome sequencing of 234 bulls facilitates mapping of monogenic and complex traits in cattle. *Nature Genetics* **46**(8), 858. doi:10.1038/ng.3034
- Daetwyler H, Xiang R, Yuan Z, Bolormaa S, Vander Jagt C, Hayes B, van der Werf J, Pryce J, Chamberlain A, Macleod I (2019) Integration of functional genomics and phenomics into genomic prediction raises its accuracy in sheep and dairy cattle. In 'Proceedings of the Association for the Advancement of Animal Breeding and Genetics', Armidale, NSW, Australia. pp. 11–14.
- de las Heras-Saldana S, Lopez BI, Moghaddar N, Park W, Park J-e, Chung KY, Lim D, Lee SH, Shin D, van der Werf JH (2020) Use of gene expression and whole-genome sequence information to improve the accuracy of genomic prediction for carcass traits in Hanwoo cattle. *Genetics, Selection, Evolution* **52**(1), 54. doi:10.1186/s12711-020-00574-2
- Erbe M, Hayes B, Matukumalli L, Goswami S, Bowman P, Reich C, Mason B, Goddard M (2012) Improving accuracy of genomic predictions within and between dairy cattle breeds with imputed high-density single nucleotide polymorphism panels. *Journal of Dairy Science* **95**(7), 4114–4129. doi:10.3168/jds.2011-5019
- Fang L, Sahana G, Ma P, Su G, Yu Y, Zhang S, Lund MS, Sørensen P (2017a) Exploring the genetic architecture and improving genomic prediction accuracy for mastitis and milk production traits in dairy cattle by mapping variants to hepatic transcriptomic regions responsive to intra-mammary infection. *Genetics, Selection, Evolution* **49**(1), 44. doi:10.1186/s12711-017-0319-0
- Fang L, Sahana G, Ma P, Su G, Yu Y, Zhang S, Lund MS, Sørensen P (2017b) Use of biological priors enhances understanding of genetic architecture and genomic prediction of complex traits within and between dairy cattle breeds. *BMC Genomics* **18**(1), 604. doi:10.1186/s12864-017-4004-z
- Fang L, Cai W, Liu S, Canela-Xandri O, Gao Y, Jiang J, Rawlik K, Li B, Schroeder SG, Rosen BD (2020) Comprehensive analyses of 723 transcriptomes enhance genetic and biological interpretations for complex traits in cattle. *Genome Research* **30**(5), 790–801. doi:10.1101/gr.250704.119
- Fink T, Lopdell TJ, Tiplady K, Handley R, Johnson TJ, Spelman RJ, Davis SR, Snell RG, Littlejohn MD (2020) A new mechanism for a familiar mutation–bovine DGAT1 K232A modulates gene expression through multi-junction exon splice enhancement. *BMC Genomics* **21**(1), 591. doi:10.1186/s12864-020-07004-z
- Fuchsberger C, Abecasis GR, Hinds DA (2015) minimac2: faster genotype imputation. *Bioinformatics* **31**(5), 782–784. doi:10.1093/bioinformatics/btu704
- Hayes BJ, Daetwyler HD (2018) 1000 Bull Genomes Project to map simple and complex genetic traits in cattle: applications and outcomes. *Annual Review of Animal Biosciences* **7**, 89–102.
- Howie B, Fuchsberger C, Stephens M, Marchini J, Abecasis GR (2012) Fast and accurate genotype imputation in genome-wide association studies through pre-phasing. *Nature Genetics* **44**(8), 955. doi:10.1038/ng.2354
- Kern C, Wang Y, Xu X, Pan Z, Halstead M, Chanthavixay G, Saelao P, Waters S, Xiang R, Chamberlain A (2021) Functional annotations of three domestic animal genomes provide vital resources for comparative and agricultural research. *Nature Communications* **12**(1), 1821. doi:10.1038/s41467-021-22100-8
- Koufariotis LT, Chen Y-PP, Stothard P, Hayes BJ (2018) Variance explained by whole genome sequence variants in coding and regulatory genome annotations for six dairy traits. *BMC Genomics* **19**(1), 237. doi:10.1186/s12864-018-4617-x
- Lee SH, Van der Werf JH (2016) MTG2: an efficient algorithm for multivariate linear mixed model analysis based on genomic information. *Bioinformatics* **32**(9), 1420–1422. doi:10.1093/bioinformatics/btw012
- Li YI, van de Geijn B, Raj A, Knowles DA, Petti AA, Golan D, Gilad Y, Pritchard JK (2016) RNA splicing is a primary link between genetic variation and disease. *Science* **352**(6285), 600–604. doi:10.1126/science.aad9417

- Liu S, Fang L, Zhou Y, Santos DJA, Xiang R, Daetwyler HD, Chamberlain AJ, Cole JB, Li CJ, Yu Y, Ma L, Zhang S, Liu GE (2019) Analyses of inter-individual variations of sperm DNA methylation and their potential implications in cattle. *BMC Genomics* **20**(1), 888. doi:[10.1186/s12864-019-6228-6](https://doi.org/10.1186/s12864-019-6228-6)
- Loh P-R, Danecek P, Palamara PF, Fuchsberger C, Reshef YA, Finucane HK, Schoenherr S, Forer L, McCarthy S, Abecasis GR (2016) Reference-based phasing using the Haplotype Reference Consortium panel. *Nature Genetics* **48**(11), 1443. doi:[10.1038/ng.3679](https://doi.org/10.1038/ng.3679)
- Lopdell TJ, Tiplady K, Struchalin M, Johnson TJ, Keehan M, Sherlock R, Couldrey C, Davis SR, Snell RG, Spelman RJ (2017) DNA and RNA-sequence based GWAS highlights membrane-transport genes as key modulators of milk lactose content. *BMC Genomics* **18**(1), 968. doi:[10.1186/s12864-017-4320-3](https://doi.org/10.1186/s12864-017-4320-3)
- MacLeod I, Bowman P, Vander Jagt C, Haile-Mariam M, Kemper K, Chamberlain A, Schrooten C, Hayes B, Goddard M (2016) Exploiting biological priors and sequence variants enhances QTL discovery and genomic prediction of complex traits. *BMC Genomics* **17**(1), 144. doi:[10.1186/s12864-016-2443-6](https://doi.org/10.1186/s12864-016-2443-6)
- McLaren W, Gil L, Hunt SE, Riat HS, Ritchie GRS, Thormann A, Flicek P, Cunningham F (2016) The Ensembl Variant Effect Predictor. *Genome Biology* **17**(1), 122. doi:[10.1186/s13059-016-0974-4](https://doi.org/10.1186/s13059-016-0974-4)
- Prowse-Wilkins CP, Wang J, Xiang R, Garner JB, Goddard ME, Chamberlain AJ (2021) Putative causal variants are enriched in annotated functional regions from six bovine tissues. *Frontiers in Genetics* **12**(1027),
- Rosen BD, Bickhart DM, Schnabel RD, Koren S, Elsik CG, Tseng E, Rowan TN, Low WY, Zimin A, Couldrey C (2020) De novo assembly of the cattle reference genome with single-molecule sequencing. *GigaScience* **9**(3), gaaa021. doi:[10.1093/gigascience/gaaa021](https://doi.org/10.1093/gigascience/gaaa021)
- Sargolzaei M, Chesnais JP, Schenkel FS (2014) A new approach for efficient genotype imputation using information from relatives. *BMC Genomics* **15**(1), 478. doi:[10.1186/1471-2164-15-478](https://doi.org/10.1186/1471-2164-15-478)
- Silva DB, Fonseca LF, Pinheiro DG, Magalhães AF, Muniz MM, Ferro JA, Baldi F, Chardulo LA, Schnabel RD, Taylor JF (2020) Spliced genes in muscle from Nelore Cattle and their association with carcass and meat quality. *Scientific Reports* **10**(1), 14701. doi:[10.1038/s41598-020-71783-4](https://doi.org/10.1038/s41598-020-71783-4)
- Weissbrod O, Hormozdiari F, Benner C, Cui R, Ulirsch J, Gazal S, Schoech AP, Van De Geijn B, Reshef Y, Márquez-Luna C (2020) Functionally informed fine-mapping and polygenic localization of complex trait heritability. *Nature Genetics* **52**(12), 1355–1363. doi:[10.1038/s41588-020-00735-5](https://doi.org/10.1038/s41588-020-00735-5)
- Xiang R, Hayes BJ, Vander Jagt CJ, MacLeod IM, Khansefid M, Bowman PJ, Yuan Z, Prowse-Wilkins CP, Reich CM, Mason BA, Garner JB, Marett LC, Chen Y, Bolormaa S, Daetwyler HD, Chamberlain AJ, Goddard ME (2018) Genome variants associated with RNA splicing variations in bovine are extensively shared between tissues. *BMC Genomics* **19**(1), 521. doi:[10.1186/s12864-018-4902-8](https://doi.org/10.1186/s12864-018-4902-8)
- Xiang R, Berg Id, MacLeod IM, Hayes BJ, Prowse-Wilkins CP, Wang M, Bolormaa S, Liu Z, Rochfort SJ, Reich CM, Mason BA, Vander Jagt CJ, Daetwyler HD, Lund MS, Chamberlain AJ, Goddard ME (2019) Quantifying the contribution of sequence variants with regulatory and evolutionary significance to 34 bovine complex traits. *Proceedings of the National Academy of Sciences of the United States of America* **116**(39), 19398–19408. doi:[10.1073/pnas.1904159116](https://doi.org/10.1073/pnas.1904159116)
- Xiang R, van den Berg I, MacLeod IM, Daetwyler HD, Goddard ME (2020) Effect direction meta-analysis of GWAS identifies extreme, prevalent and shared pleiotropy in a large mammal. *Communications Biology* **3**(1), 88. doi:[10.1038/s42003-020-0823-6](https://doi.org/10.1038/s42003-020-0823-6)
- Xiang R, MacLeod IM, Daetwyler HD, de Jong G, O'Connor E, Schrooten C, Chamberlain AJ, Goddard ME (2021) Genome-wide fine-mapping identifies pleiotropic and functional variants that predict many traits across global cattle populations. *Nature Communications* **12**(1), 860. doi:[10.1038/s41467-021-21001-0](https://doi.org/10.1038/s41467-021-21001-0)
- Xu L, Gao N, Wang Z, Xu L, Liu Y, Chen Y, Xu L, Gao X, Zhang L, Gao H (2020) Incorporating genome annotation into genomic prediction for carcass traits in Chinese simmental beef cattle. *Frontiers in Genetics* **11**, 481.

Handling editor: Sue Hatcher