

BayesR3 enables fast MCMC blocked processing for largescale multi-trait genomic prediction and QTN mapping analysis

Edmond J. Breen¹, Iona M. MacLeod¹, Phuong N. Ho¹, Mekonnen Haile-Mariam¹, Jennie E. Pryce^{1,2}, Carl D. Thomas¹, Hans D. Daetwyler^{1,2}, and Michael E. Goddard^{1,3}

¹ Agriculture Victoria, AgriBio, Centre for AgriBioscience, Bundoora, Victoria 3083, Australia

² School of Applied Systems Biology, La Trobe University, Bundoora, Victoria 3083, Australia

³ Faculty of Veterinary and Agricultural Sciences, University of Melbourne, Parkville, VIC, 3052 Australia

Supplementary Note 1

Mixture Model Distribution

The mathematics used for BayesR3 is given in detail in the methods section of the main manuscript, but briefly the SNP effects are modelled by a mixture of four normal distributions with zero mean and increasing variances as specified by:

$$p(g_j | \pi, \sigma_g^2) = \pi_1 \times N(0, 0 \times \sigma_g^2) + \pi_2 \times N(0, 10^{-4} \times \sigma_g^2) + \pi_3 \times N(0, 10^{-3} \times \sigma_g^2) + \pi_4 \times N(0, 10^{-2} \times \sigma_g^2) \quad (1)$$

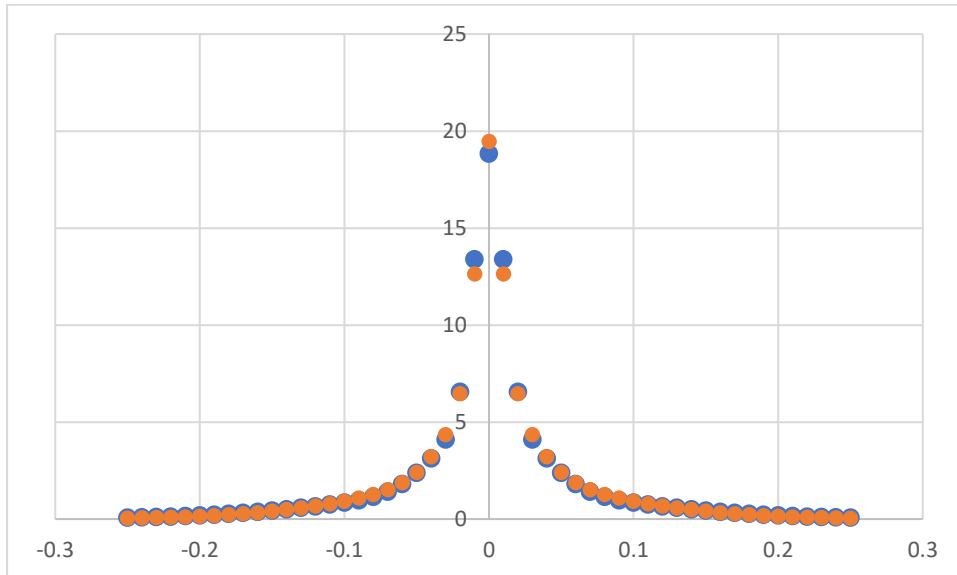
Where σ_g^2 is the additive genetic variance explained by the SNPs cumulatively and is estimated from the data. The mixing proportions π are also estimated from the data and are assumed to be drawn from a Dirichlet distribution with parameter = (1,1,1,1), a uniform prior, such that any SNP a priori is equally likely to be assigned to any one of the 4 distributions. The choice of 4 distributions, is historical (1), but any number of distributions can be used. For example, in very large datasets, adding the variance group $10^{-5} \times \sigma_g^2$ can help capture SNPs with very small effects (2). However, the allocations values (0, 10^{-4} , 10^{-3} , 10^{-2}) seen in Equation (1) can mimic a broad range of parametric distributions, such as a t or a reflected gamma, by varying the proportions π in each distribution. They can describe a distribution with long tails as we expect for SNP effects where there are many small effects and the occasional large effect (Figure S1).

The 10x scaling between the allocation values is arbitrary but convenient in practice. It allows the distributions generated to be relatively smooth and effects can shuffle from one distribution to the next between MCMC cycles. Figure S1 shows a distribution that is a mixture of 3 normal distributions with variances (0.0001, 0.001 and 0.01) in blue compared to a mixture of 6 distribution with variances 0.00005, 0.0001, 0.0005, 0.001, 0.005, 0.01 in red. They are very similar and the use of either one as the prior would have little effect on the resulting estimated SNP effects.

These allocation variances could be estimated from the data, but this is unnecessary and introduces additional complexity. The fact that similar distributions can be generated by different mixtures is a warning that the data cannot distinguish a variety of possible distributions because they are essentially the same. Also, if the variances were sampled within the MCMC process the large variance and small variance would at times swop, otherwise the chain is not mixing fully. This makes it difficult to interpret the summary statistics from the chain.

38 Therefore, we think that allowing the proportions π and σ_g^2 to be derived from the data gives the
39 model ample flexibility to fit a variety of useful distributions. It is also worth noting that Bayes R has
40 been published many times with this arrangement of variances.

41



42

43 *Figure S1: Mixture distributions. A distribution made up of equal parts of 0.0001, 0.001 and 0.01 variances (blue dots) and a*
44 *2nd mixture distribution made from 6 normal distributions with variances 0.00005, 0.0001, 0.0005, 0.001, 0.005, 0.01 in*
45 *proportions 0.11, 0.16, 0.16, 0.17, 0.17, 0.23 in orange.*

46 Supplementary References

- 47 1. Erbe M, Hayes BJ, Matukumalli LK, Goswami S, Bowman PJ, Reich CM, et al. Improving
48 accuracy of genomic predictions within and between dairy cattle breeds with imputed high-density
49 single nucleotide polymorphism panels. *J Dairy Sci.* 2012;95(7):4114-29.
- 50 2. Moser G, Lee SH, Hayes BJ, Goddard ME, Wray NR, Visscher PM. Simultaneous discovery,
51 estimation and prediction analysis of complex traits using a bayesian mixture model. *PLoS Genet.*
52 2015;11(4):e1004969.

53

54 Supplementary Note 2

55 Determining an Optimal Block Size

56 Here we look at two methods for determining the optimal block size and one method to determining
57 the number of inner iterations.

58

59 Optimizing Block Size and the Number of Inner Cycles by Constraining Accuracy and 60 Minimizing Time.

61 In determining the optimal block size, it is noted with respect to a given chain length, that as block
62 size increases BayesR3 prediction accuracy decreases, albeit only slightly, (see Figure 4 main
63 manuscript), as processing time rapidly decreases. This occurs, until a certain block size is reached,
64 after which processing time remains essentially constant and at times may even increase as block
65 size increases. Therefore, we looked for a simple method to define in generally terms the optimal
66 block size and number of inner iterations.

67 Let

68 $n_M =$ No. of SNPs
69 $n =$ No. of SNPs per block
70 $\frac{n_M}{n} =$ No. of blocks
71 $n_R =$ No. of Records
72 $x =$ No. of inner cycles
73 $y =$ No. of outer cycles
74 $xy =$ total no. of cycles

75 The number of blocks is determined by the smallest integer larger or equal to the ratio $\frac{n_M}{n}$, this
76 means that all blocks will be the same size except for the last block, which can be smaller.

78

77 **Time taken:**

79 The major contributors to the processing time are:

- 80 1. Calculating the sum for each SNP when processing a new block. This is proportional to:
81 $yn_M n_R$.
82 2. Cycling around a block, sampling a SNP effect, and updating the total for all other SNPs in
83 the block. This is proportional to $yn_M xn$

84 Then:

85
$$\text{Time} = yn_M n_R + yn_M xn$$

86

87 **Accuracy:**

88 We wish to take enough samples of each SNP effect to reduce the Monte Carlo sampling error of the
89 mean to an acceptably low value. Each individual sampled value of a SNP effect (\hat{b}) can be modelled
90 as:

91
$$\hat{b}_{ij} = b + u_i + e_{ij}$$

92

93 where \hat{b}_{ij} is the sampled value of b in outer cycle i and inner cycle j . $u_i =$ effect common to all
94 samples in outer cycle i , and e_{ij} is the effect of the inner cycle j within outer cycle i .

95

96 The variation of the samples within the same block are correlated because they are based on the
97 sampled values from all other blocks. This correlation is the reason why u_i is included in the above
98 formula for \hat{b}_{ij} . The importance of u_i depends on the LD between SNPs in the current block and
99 SNPs in all other blocks. This LD decreases as block size increases. Therefore, we will assume:

100

101
$$v(u_i) \propto \frac{1}{n}$$

102
$$v(\hat{b}_{ij}) = \frac{\sigma^2}{n} + \sigma^2$$

103

104 and the Monte Carlo variance of the average \hat{b}_{ij} is: $\frac{\sigma^2}{yn} + \frac{\sigma^2}{xy}$.

105

106 We wish to optimize x and n by minimizing the time taken while holding constant the Monte Carlo
 107 sampling variance. Using a Lagrange multiplier, the objective is:
 108

$$109 \quad yn_M(n_R + xn) + \lambda\left(\frac{1}{ny} + \frac{1}{xy} - k\right)$$

$$110 \quad \frac{\delta}{\delta n} = yn_Mx - \lambda \frac{1}{n^2y} = 0 \Rightarrow n^2 = \frac{\lambda}{n_Mxy^2}$$

$$111 \quad \frac{\delta}{\delta x} = yn_Mn - \lambda \frac{1}{x^2y} = 0 \Rightarrow x^2 = \frac{\lambda}{n_Mny^2}$$

112
 113 Therefore, $n = x$ and the constraint becomes $\frac{2}{ny} - k = 0$, and hence

$$114$$

$$115 \quad y = \frac{2}{nk}$$

$$116 \quad \text{Time} = yn_M(n_R + xn)$$

$$117 \quad = \frac{2n_M}{nk}(n_R + n^2)$$

$$118 \quad = \frac{2n_M}{k}\left(\frac{n_R}{n} + n\right)$$

$$119 \quad \frac{\delta}{\delta n} = \frac{2n_M}{k}\left(\frac{-n_R}{n^2} + 1\right) = 0$$

$$120 \quad \Rightarrow n = \sqrt{n_R}$$

121
 122 Thus, the approximate optimum is $n = x = \sqrt{n_R}$. This is approximate because of the assumptions
 123 made and the exact optimum may be data set dependent. However, it is intuitively reasonable: If
 124 there are many SNPs per block, the value for each SNP depends on the current value if all the other
 125 SNPs in the block so it is worthwhile to take many cycles around the block. The time saved by
 126 blocking is due to processing the individual records only once per block and therefore as n_R
 127 increases the optimum n also increases.
 128

129 [Determining the optimal block size using a curvature equation](#)

130 Given the number of inner iterations is set to equal the number of SNPs within a block, we note the
 131 observed change in processing time with respect to block size for a given genomic data set as shown
 132 in Figure 5b (main manuscript), was successfully model using the function $f(n) = \frac{n_R+n}{n}$. This
 133 function has one point on the curve where the curvature is a maximum, and which corresponds to a
 134 transition from high to low curvature. The derivative of $f(n)$ is $f'(n) = \frac{-n_R}{n^2}$, $f''(n) = \frac{2n_R}{n^3}$, and the
 135 curvature of the curve, in Figure 5c main manuscript, is given by $\kappa(n) = \left| \frac{f''(n)}{(1+(f'(n))^2)^{\frac{3}{2}}} \right|$. When the
 136 derivative of κ is set to zero, it has a corresponding positive root where the curvature is maximized
 137 and where $n = \sqrt{n_R}$. In terms of optimization this represents an elbow or knee point and is the
 138 optimum between the benefit achieved between a reduced processing time and the loss in

139 prediction accuracy as n increases. Therefore, we suggest that the block size should not increase
 140 beyond $\sqrt{n_R}$.
 141

142 **Supplementary Tables**

143

144 *Supplementary Table 1: Details of the QTL annotated in **Error! Reference source not found.** for milk composition traits*
 145 *discovered in the multi-trait BayesR3 MIR analysis as well as the multi-trait Milk, Fat and Protein Yield BayesR3 (MFP_BR3)*
 146 *and BayesR3C (MFP_BR3C) analyses. Details include the underlying candidate genes and previously reported overlapping*
 147 *milk trait QTL.*

QTL midpoint position (bp: see Error! Reference source not found.)	Multi-trait analysis detecting the QTL	Candidate gene(s) and start...stop position (bp)	Examples of published reports for milk traits	Examples of milk traits previously reported for this QTL position
Chr1:142827704	MIR	<i>SLC37A1</i> 142772300...142873917	(30, 41)	P, Milk Yield P, Mg, K, Na
Chr2:131204809	MIR	<i>ALPL</i> 131181421..131268191	(30)	Na
Chr3:15387272-15484820	MIR MFP_BR3 MFP_BR3C	<i>SLC50A1</i> 15463149...15465593 <i>DPM3</i> 15462231...15462806 <i>EFNA1</i> 15466565...15473164 <i>LOC107132270</i> (ncRNA) 15465686...15496399	(30, 31)	Lactose percentage Na
Chr5:93534138-93538860	MIR MFP_BR3 MFP_BR3C	<i>MGST1</i> 93495438...93520998 <i>SLC15A5</i> 93602194...93699207	(25, 31, 42)	fat percentage Lactose yield Milk and Fat yield, Fat percentage
Chr6:45052030	MIR	<i>LOC112447058</i> ncRNA 45060561...45065030	(30)	P, K
Chr6:85419916-85475175	MIR MFP_BR3 MFP_BR3C	<i>CSN2</i> 85449173...85457867 <i>CSN1S1</i> 85411601...85429256	(25, 30, 41)	protein yield, protein percent Na, Ca Milk and Protein Yield
Chr11:103243985	MIR MFP_BR3 MFP_BR3C	<i>PAEP</i> 103255963...103260862	(30, 41, 43)	Milk, fat and Protein Yield, Fat percentage Protein percentage Citrate, P
Chr11:104186080	MFP_BR3 MFP_BR3C	<i>ABO</i> 104176830...104214758	(25, 44)	Protein yield oligosaccharides
Chr14: 263754 - 792410	MIR MFP_BR3 MFP_BR3C	<i>DGAT1</i> 603981...614153	(29, 45)	Milk, fat and protein yield, fat and protein percentage

		(many other genes including solute carriers <i>SLC39A4</i> and <i>SLC52A2</i>)		
Chr15:81098284	MFP_BR3 MFP_BR3C	<i>CTNND1</i> and olfactory receptor genes	No reports found	-
Chr16:1750054 - 1803090	MFP_BR3 MFP_BR3C	<i>SOX13</i> 1867662...1912526 <i>LOC104974354</i> 1791622...1806028	No reports found	-
Chr20:31887560	MFP_BR3 MFP_BR3C	<i>GHR</i> 31869704...32043372	(43, 46, 47)	Milk, Fat and Protein Yield, Fat and Protein percentage, Milk Yield Protein percent
Chr20:58389561	MIR MFP_BR3 MFP_BR3C	<i>ANKH</i> 58307527...58477499	(30, 31, 41)	Fat yield, protein percentage Lactose percentage K
Chr27:36499460	MIR	<i>GPAT4 (AGPAT6)</i> 36508780...36539760	(30, 48)	Fat and protein percentage, protein yield, lactose yield and percentage, milk fatty acids Mg

148

149

150

Supplementary Table 2: Distribution of SNP effects on PCs derived from milk, fat, and protein yield and PC heritabilities.

151

Note the loadings give the coefficients of the linear combinations of the centered and scaled continuous variables for milk,

152

fat, and protein yield. Distributions 1-4 are each normal distributions with mean = 0 and variances = 0, 0.0001, 0.001, 0.01

153

times the genetic variance, respectively.

PC	PC loadings			Distribution				h^2
	Milk	Fat	Protein	k_1	k_2	k_3	k_4	
1	-0.603	-0.493	0.627	140	9801	5	2	0.42
2	-0.443	0.861	0.251	2688	7219	12	30	0.48
3	-0.663	-0.126	-0.738	2293	7555	79	21	0.52

154

155 *Supplementary Table 3: Distribution of SNP effects on PCs derived from milk, fat and protein yield when using prior*
 156 *information from analysis of milk MIR spectra.*

Class	Total Number of SNPs	Number of SNPs in Model	PC	Distribution			
				k_1	k_2	k_3	k_4
1	992	48 (4.8%)	1	7	30	9	2
			2	5	4	27	11
			3	2	2	32	12
2	238677	4084 (1.6%)	1	138	3940	5	1
			2	614	3460	8	12
			3	348	3660	58	5
3	238667	4061 (1.6%)	1	438	3520	3	0
			2	1880	2100	4	4
			3	1340	2640	5	3
4	157535	2665 (1.5%)	1	202	2400	3	0
			2	1140	1440	2	0
			3	1560	1036	2	0

157

158

159 *Supplementary Table 4: Multi-Trait Analysis of Milk Production Traits of Dairy Cattle – accuracy of prediction using BayesR3*
 160 *and BayesR3C. BayesR3C used the multi-trait MIR Q2 probabilities to allocate variants to four classes (see Materials &*
 161 *Methods in full paper). Reference N=65,637, & Validation populations were as described for the Single Trait Analyses:*
 162 *HOL_Bull was 702 Holstein bulls, JER_Bull was 675 Jersey bulls and RDC_Cows included 3082 Australian Red cows. Accuracy*
 163 *was averaged across 5 MCMC chains for each PC trait.*

Validation Set	PC Trait	BayesR3 Accuracy	BayesR3C Accuracy
Hol_Bull	1	0.708	0.710
	2	0.819	0.818
	3	0.841	0.840
JER_Bull	1	0.753	0.753
	2	0.799	0.800
	3	0.852	0.854
RDC_Cows	1	0.235	0.231
	2	0.396	0.384
	3	0.255	0.263

164

165