

Impact of rare variants on the quality of genomic prediction in dairy cattle

Tomasz Suchocki and Joanna Szyda

Wroclaw University of Environmental and Life Sciences
Department of Genetics
Biostatistics Group



Motivation:

- ▶ SNP selection based on minor allele frequencies and call rate.
- ▶ What with SNPs with very good call rate and very low minor allele frequencies? Will they influence on accuracy of genomic breeding values?

Aim of the study is:

- ▶ identification SNP markers with rare allelic variants;
- ▶ comparison of accuracy of genomic breeding values for data sets with and without rare allelic variants;
- ▶ comparison of Interbull validation procedure for data sets with and without rare allelic variants.

Animals:

- ▶ 5 068 individuals
 - ▶ 3 100 proven bulls
 - ▶ 1 968 young bulls

Traits:

- ▶ Production: milk, fat and protein yield
- ▶ Fertility: non-return rate of heifers and cows
- ▶ Udder health: somatic cell score

Genotype:

- ▶ 46 267 SNPs after selection based on $MAF > 1\%$ and call rate $> 95\%$
- ▶ 53 862 SNPs without MAF selection

SNP effect estimation:

$$y = \mu + Zq + \epsilon,$$

- ▶ y - deregressed EBV
- ▶ μ - overall mean
- ▶ q - random SNP effect $\sim \mathcal{N}(0, I \cdot \frac{\hat{\sigma}_\alpha^2}{N_{SNP}})$
- ▶ $N_{SNP} = 46\,267$ or $N_{SNP} = 53\,862$
- ▶ $Z = \{-1, 0, 1\}$
- ▶ ϵ - error term $\sim \mathcal{N}(0, D \cdot \hat{\sigma}_\epsilon^2)$

SNP effect estimation:

$$y = \mu + Zq + \epsilon,$$

- ▶ y - deregressed EBV
- ▶ μ - overall mean
- ▶ q - random SNP effect $\sim \mathcal{N}(0, I \cdot \frac{\hat{\sigma}_\alpha^2}{N_{SNP}})$
- ▶ $N_{SNP} = 46\,267$ or $N_{SNP} = 53\,862$
- ▶ $Z = \{-1, 0, 1\}$
- ▶ ϵ - error term $\sim \mathcal{N}(0, D \cdot \hat{\sigma}_\epsilon^2)$

$$DGV = Z \cdot \hat{q}$$

Calculation of reliability:

$$Rel = diag \left\{ \left(Q - \frac{\hat{\sigma}_{\epsilon}^2}{\hat{\sigma}_{\alpha}^2} C^{22} \right) Q^{-1} \right\},$$

- ▶ C^{22} - inverse of coefficient matrix for MME
- ▶ $Q = ZZ^T \frac{1}{p_{het}^b}$
- ▶ p_{het}^b - sum over all SNP of heterozygous genotype frequencies in base population

The bias in the national genomic evaluations will be tested using a regression model:

$$\phi_i = b_0 + b_1 GEBV_i + \epsilon_i,$$

where

- ▶ ϕ_i is de-regressed predicted genetic merits or daughter deviations from the bulls that have EDC > 20
- ▶ *GEBV* - parent averages plus genomic prediction equations

SOURCE: <http://interbull.org>

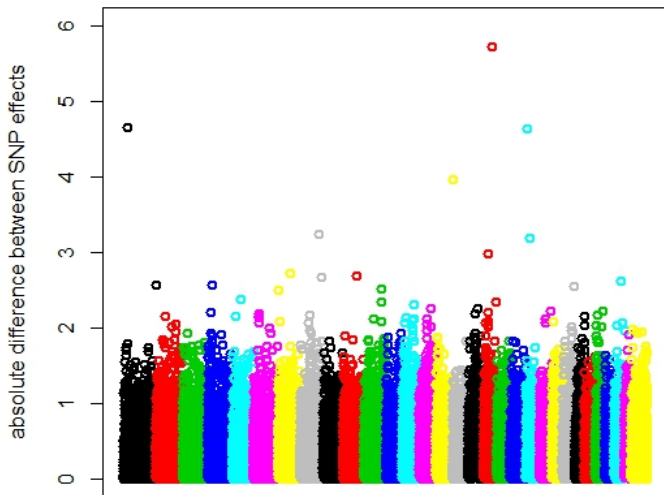
The improvement of the added genomic information to the parental information will be tested using following model:

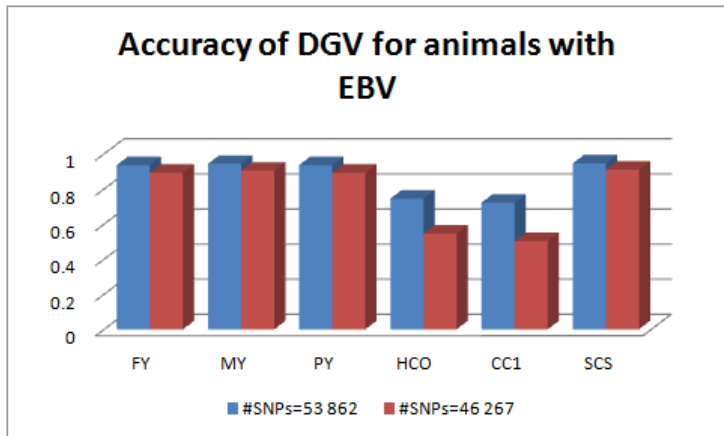
$$\phi_i = b_0 + b_1 EBV_i + \epsilon_i,$$

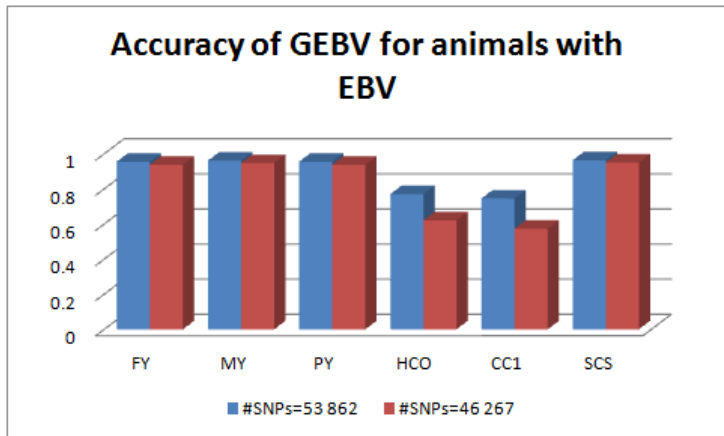
- ▶ *EBV* - genetic merit estimates based only on parent averages

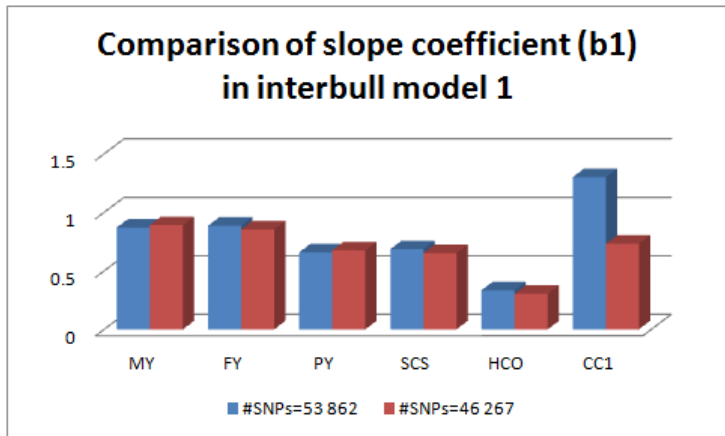
SOURCE: <http://interbull.org>

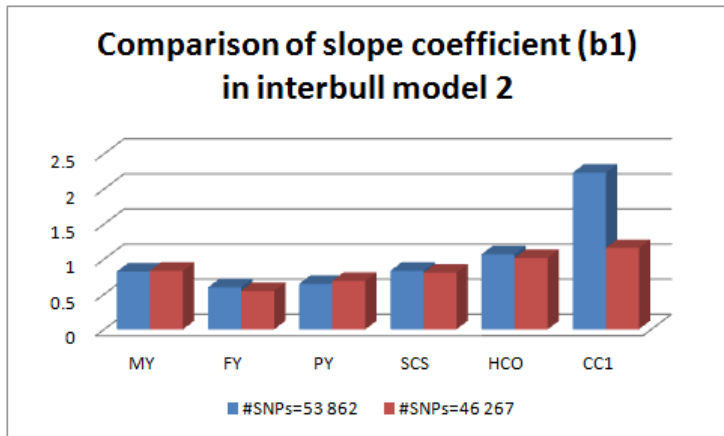
Results - Comparison of common SNPs

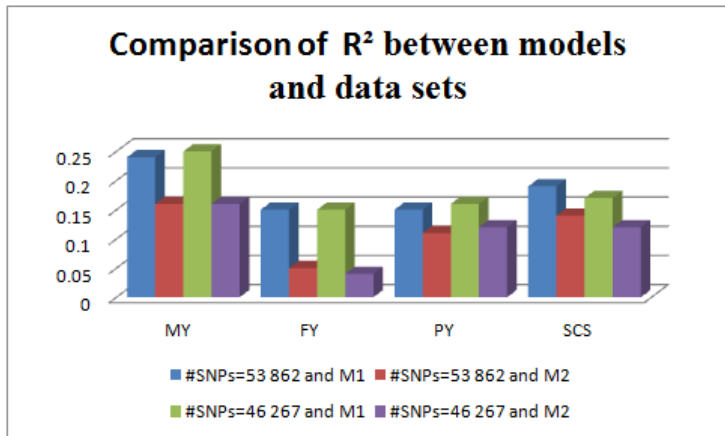












- ▶ Accuracy of GEBV and DGV is higher for data set with 53 862 SNPs.
- ▶ Especially for fertility traits increase of accuracy is high. For CC1 it is 29% for GEBV and 43% for DGV.
- ▶ For fat yield, somatic cell score and HCO the parameter b_1 from Interbull validation test is closer to $E(b_1)$ for data set with 53 862 SNPs.
- ▶ For other traits better results are for data set with 46 267 SNPs.



Thanks for MASinBULL!

We thank the members of the MASinBULL consortium, who provided the data set used in the analysis.

MASinBULL



Thank you for your attention!

