

# Performance of Genomic Prediction using Different Multi-Allelic Genomic Relationship Matrices for Genotyping-by-Sequencing (GBS) Data

Jie Kang<sup>1</sup>\*, Philip Wimmer<sup>2</sup>, Michael Steiner<sup>3</sup>, Stephen Symon<sup>4</sup>, Don Milbourne<sup>4</sup>, Marty Parfitt<sup>4</sup>, Joanne Jacobs<sup>4</sup>, Ken Dodds<sup>4</sup>

<sup>1</sup>Department of Mathematics and Statistics, University of Otago, Dunedin, New Zealand

<sup>2</sup>Department of Microbiology, University of Otago, Dunedin, New Zealand

<sup>3</sup>Tangaroa, Craig Research Centre, GSB Park, Christchurch, New Zealand

<sup>4</sup>Agriculture, Invermay Agricultural Centre, Invermay, New Zealand

## Introduction

- Genomic prediction (GP) has been frequently implemented using SNP data since its introduction almost a decade ago.
- Genomic relationship matrices (GRM) can be derived from SNP data, and used with genomic best linear unbiased prediction (GBLUP) to perform GP.
- Short haplotypes ('ShortSteps') refer to multiple variants situated in small genomic segments such as those captured within SNP reads.
- ShortSteps are more advantageous than single nucleotide polymorphisms (SNPs) because they are more likely to be strongly associated with causal variants. Multi-allelic genomic relationship matrices (MA-GRM) constructed using shortSteps are expected to be superior.

## Aim

Investigate whether performance of genomic prediction can be further improved by using multi-allelic genomic relationship matrices (MA-GRM).

## Method

Three different approaches for developing MA-GRM were assessed, alongside pedigree (A) and SNP genotype (SNP) methods, using a simulated dataset based on SNP and phenotype data from parental ryegrass.

### MA-GRM

#### 1. Haplotype-coverage method (Wimmer et al., 2014)

$$MA_{GRM} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

where  $n$  denotes the number of ordered haplotype reads (coverage), and  $x$  denotes the haplotype allele frequencies.

#### 2. Haplotype similarity method (Parfitt et al., 2016)

$$MA_{GRM} = \frac{1}{n} \sum_{i=1}^n x_i^2$$

where  $K$  is the Kronecker product between an  $n \times n$  identity matrix ( $n =$  number of individuals) and a  $1 \times 2$  unit matrix  $I$  in the genomic relationship matrix.

#### 3. Linkage disequilibrium method (Hetherington et al., 2016)

$$MA_{GRM} = \frac{1}{n} \sum_{i=1}^n x_i^2$$

where  $D$  is a block-diagonal matrix containing the (average) multiallelic variance-covariance matrix of haplotypes.

### MA-GBLUP

$$y = X\beta + Z\gamma + \epsilon$$

where  $y$  is a matrix of phenotypes,  $X$  and  $Z$  are incidence matrices,  $\beta$  and  $\gamma$  are fixed and random effects with  $\text{Var}(\gamma) = G\sigma^2$  and  $\epsilon$  is the error term. In particular,  $G$  can be any relationship matrix, including those listed above.

## References

- Wimmer et al. (2014). *Genetics Selection Evolution*, 46, 1-6.  
Parfitt et al. (2016). *Genetics Selection Evolution*, 48, 75.  
Hetherington et al. (2016). *Genetics*, 136, 100-105.  
Sims et al. (2016). *Proc 11th World Congr Genet Appl to Livest Prod*, 507.  
Parfitt et al. (2016). *Theoretical and Applied Genetics*, 128, 109-126.

## Simulation

- April 2000 data of 50 ryegrass samples (with a known pedigree, Fig. 1) were simulated using *simuPOP* (Zeng et al., 2016), based on the ryegrass reference genome described in Parfitt et al. (2016).



Figure 1: Pedigree used in simulation

- 446,100 out of 1,805,744 virtually digested SNP fragments were retained after fragment characterization (50, 100 bp, 200,507 tags contained SNPs (2 SNPs per tag on average) with 876,000 SNPs in total).
- 1000 QTL from a quantitative trait with heritability of 0.4, 0.4, 0.2, 0.8 were simulated (1000 replicates of each scenario), with the QTL effects and residuals drawn from standard normal distributions.

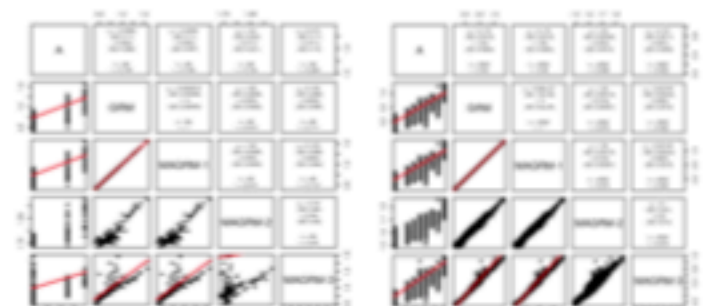


Figure 2: Scatter comparison of diagonal (diag) of different GRM relative to relationship matrix

## Results & Discussion

- The pedigree (A), SNP genotype (SNP) and three haplotype-based relationship matrices (MA-GRM) returned consistent results (Fig. 2), where off-diagonal and diagonals of all relationship matrices are highly correlated (apart from the diagonals of MA-GRM-3).
- Performance of GBLUP using different GRM was almost identical (Fig. 2). However, the results may be confounded by the population size and make-up (including SNP density and linkage disequilibrium between SNPs).
- This study assessed SNP data low linkage depth - the impacts of low sequencing depth and other SNP-specific sources of variation need to be considered in the future.

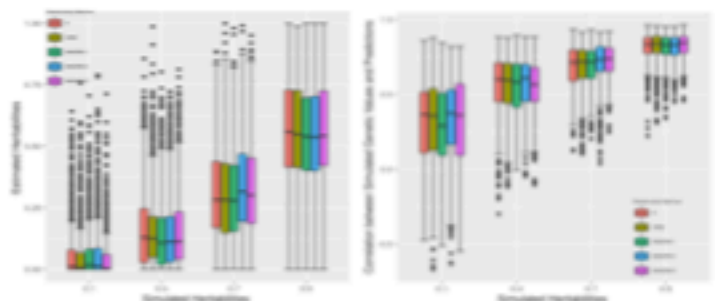


Figure 3: Performance of GBLUP using different GRM methods. The plots show the relationship between the diagonal elements of the GRM matrices and the corresponding relationship matrix.

## Conclusion

- There was a strong correlation between MA-GRM and other relationship matrices.
- The use of MA-GRM did not improve performance of GP under current settings.