

Efficient computation of base generation allele frequencies

11 February; Interbull meeting, Auckland, New Zealand

Michael Aldridge, Jeremie Vandenplas & **Mario Calus**



Allele frequencies in genomic prediction

- Genomic prediction requires **allele frequencies (AF)**
- Commonly, AF are current **data averages**
- **Theoretically**, AF should be computed for the **base generation**

Base generation AF

Base generation = base generation in pedigree!

Base generation AF required for calculation of:

- Genomic **relationships** in (single-step) GBLUP
- Model-based **reliabilities** for multi-step genomic evaluations
- Computation of **relationships** among **metafounders**¹

Objective

Compare **accuracy** and **efficiency**
of different methods to compute
base generation **allele frequencies**

Methods – overview

- AF: $p = \frac{1}{2}\hat{\mu}$

Method	Mean is estimated:
All	Across all genotypes
Oldest	Across oldest generation genotyped
BLUP	In BLUP model
GLS	General Least Squares (GLS)

Methods - BLUP

- BLUP model; $y = \text{genotype } (0,1,2)$
- $h^2=0.99$; allowing some genotyping error
- Univariate; or multivariate with **zero** genetic correlations
- Implemented using MiXBLUP

Methods – GLS (dense / sparse)

- GLS: $\hat{\mu}_i = (\mathbf{1}'\mathbf{A}_{22}^{-1}\mathbf{1})^{-1}\mathbf{1}'\mathbf{A}_{22}^{-1}\mathbf{Z}_i$

- **Dense:** Compute and invert \mathbf{A}_{22}

Calc_grm

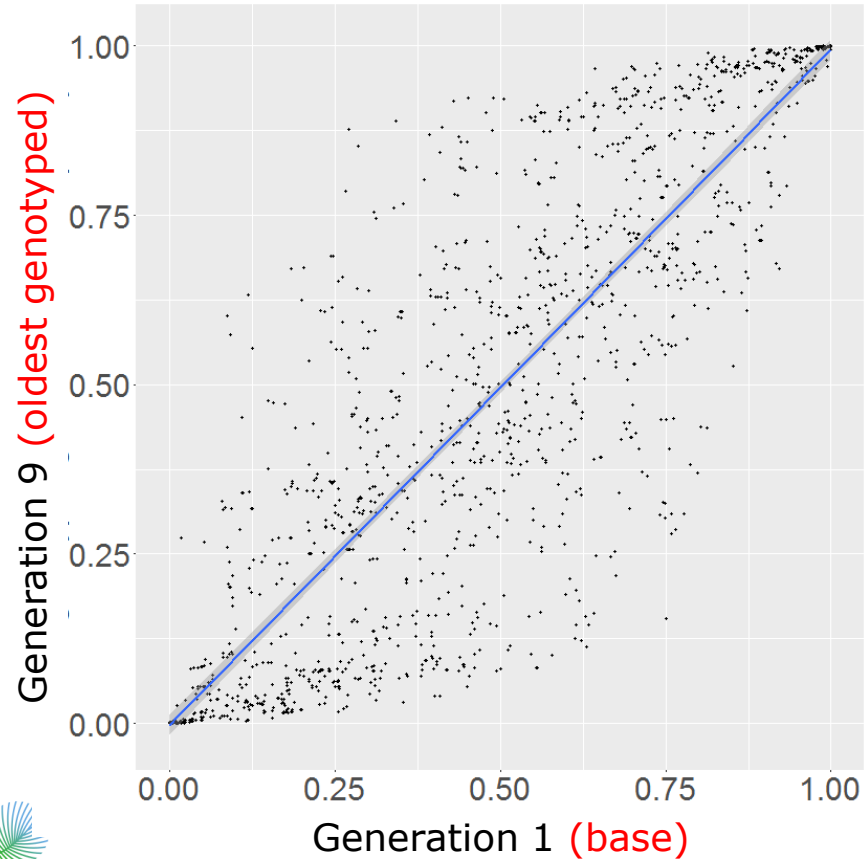
- **Sparse:** $\mathbf{A}_{22}^{-1}\mathbf{1} = \left(\mathbf{A}^{22} - \mathbf{A}^{21}(\mathbf{A}^{11})^{-1}\mathbf{A}^{12}\right)\mathbf{1}$

Own program / Intel MKL-PARDISO

Data (simulation)

- Holstein-like population
- Generations 9 to 12 (after base) fully genotyped
- 325,266 animals in pedigree; 100,078 genotyped
- 1670 SNPs (providing replication)
- Selection: **None** or **Strong**

Change in AF across generations (with selection)



Results - accuracy

Method	Without selection	With selection
All	0.99 ± 0.01	0.87 ± 0.01
Oldest	0.99 ± 0.01	0.88 ± 0.01
BLUP	0.99 ± 0.01	0.96 ± 0.01
GLS_dense	0.99 ± 0.01	0.97 ± 0.01
GLS_sparse	0.99 ± 0.01	0.97 ± 0.01

Results - efficiency

Method	Process time	RAM
All	0-00:03:44	7.8 GB
Oldest	0-00:01:19	1.6 GB
BLUP (60 SNPs)	0-13:42:17	49.0 GB
GLS_dense	50-20:12:16	165.9 GB
GLS_sparse	0-00:01:28	2.6 GB

=> Efficiency of GLS_sparse is very competitive!

Discussion

- Few GLS_sparse estimates outside 0-1 range:
 - Only for **very low MAF** <0.001
 - **Swapping** allele code solved most of those

- Estimates were not affected when having:
 - 2% genotyping **errors**
 - 25% of sires **unknown**

Conclusions

- Base generation AF required for:
 - Genomic **relationships** in (single-step) GBLUP
 - Model-based **reliabilities** for multi-step genomic evaluations
 - Computation of **relationships** among **metafounders**

- **GLS_sparse** estimator recommended
 - Accurate & very efficient

Acknowledgements

