

# A single-step principal component ridge regression model for large-scale genomic evaluations

J. Ødegård<sup>1,2</sup> and T. H. E. Meuwissen<sup>2</sup>

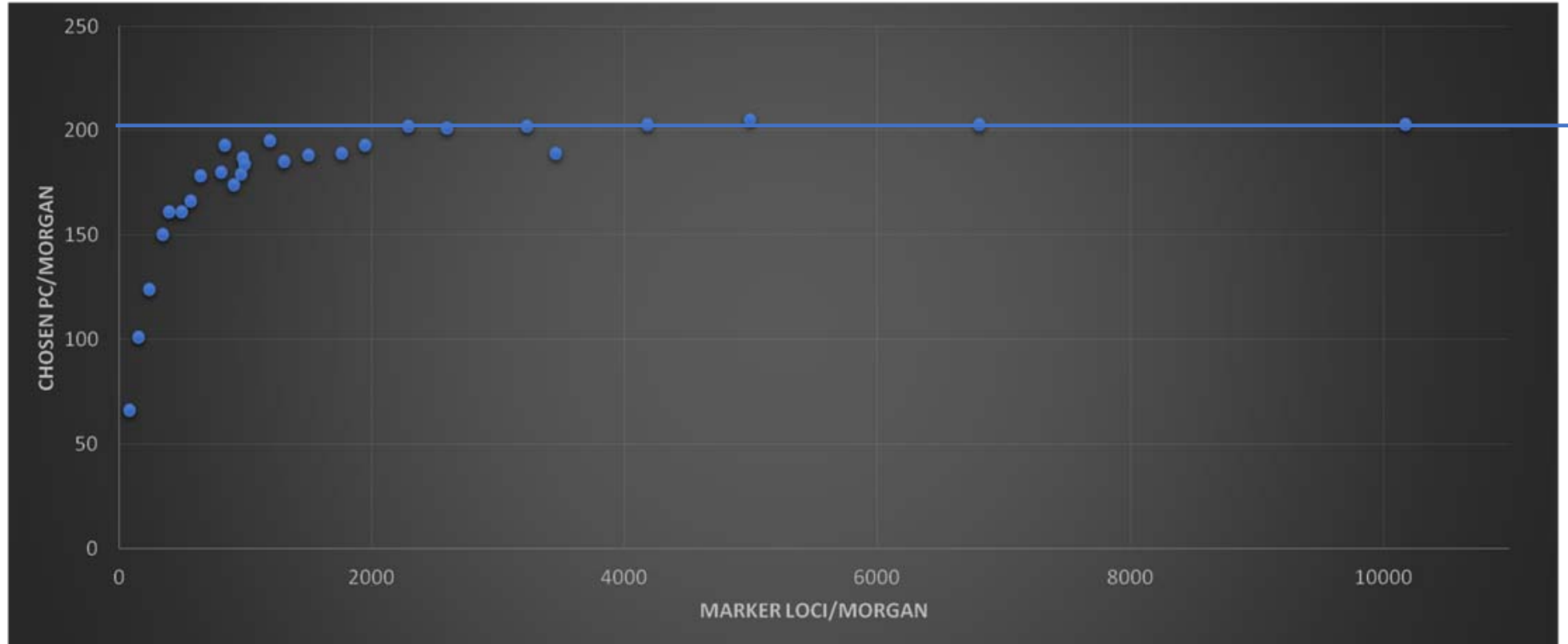
<sup>1</sup> AquaGen AS, P.O. Box 1240, N-7462 Trondheim, Norway

<sup>2</sup> Norwegian University of Life Sciences, P.O. Box 5003, N-1432 Aas, Norway

# Large-scale genomic data

- Original single-step model (ssGBLUP), Legarra et al. (2009), Christensen & Lund (2010)
  - Complexity increases with number of genotyped animals
  - Inverse genomic relationship matrix (GRM) must be computed prior to the analysis
- Single-step marker effects model (ssMEM), Fernando et al. (2016)
  - No need for inverse GRM
  - Complexity depends on number of loci
- Populations of limited  $N_e$ 
  - Limited number of haplotypes
- Genomic data can be approximated by a smaller number of principal components

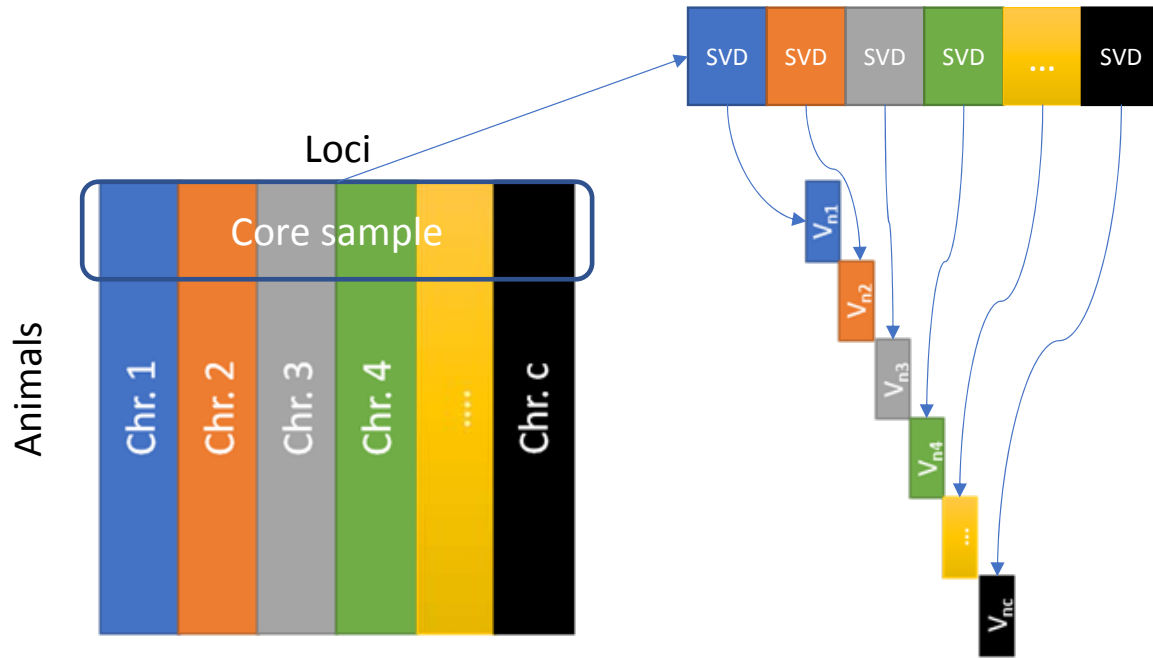
Principal components explaining >99% of variance  
( $N_e = 500$ ,  $N = 10,000$ )



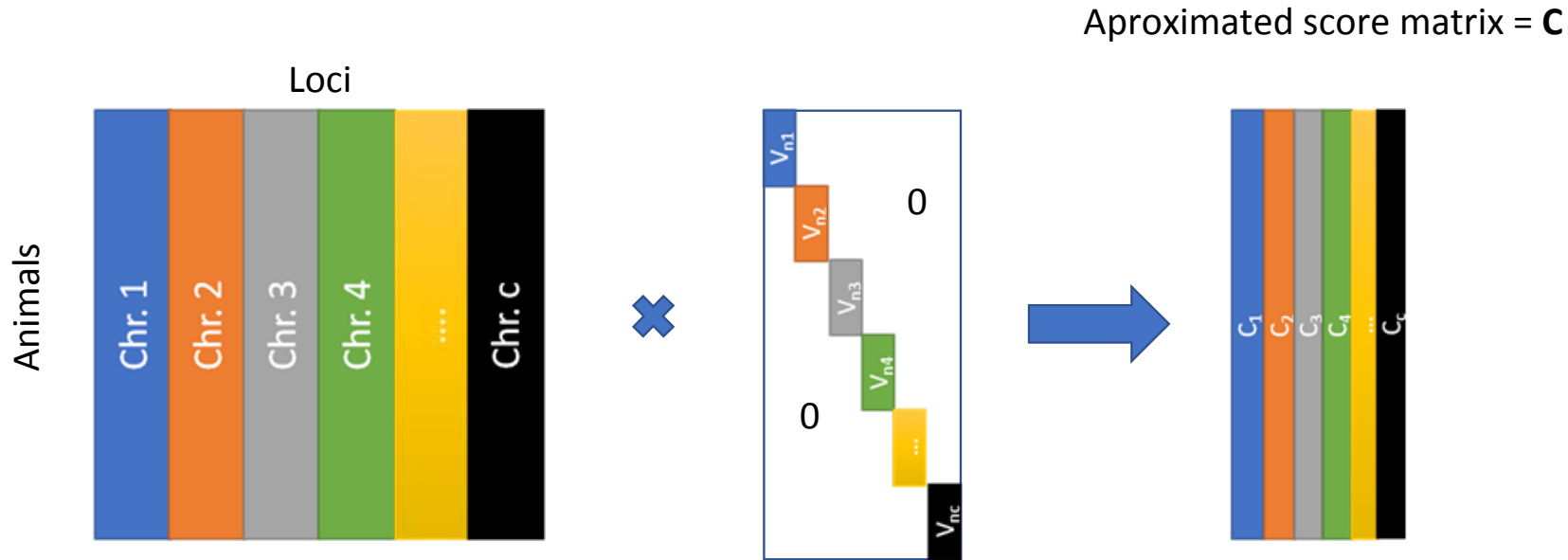
# Singular value decomposition (SVD) of genomic data

- SVD of  $N \times k$  (centered) genotype matrix
  - $\mathbf{M} = \mathbf{U}\mathbf{S}\mathbf{V}'$
  - $\mathbf{U}$  =eigenvectors of  $\mathbf{M}\mathbf{M}'$  (orthonormal),  $\mathbf{U}'\mathbf{U} = \mathbf{I}$
  - $\mathbf{V}$  =eigenvectors of  $\mathbf{M}'\mathbf{M}$ (orthonormal),  $\mathbf{V}'\mathbf{V} = \mathbf{I}$
  - $\mathbf{S}$  is a diagonal matrix (square root of eigenvalues)
- Principal component ridge regression model
  - $\mathbf{y} = \mathbf{T}\mathbf{s} + \mathbf{e}$ 
    - $\mathbf{s} = \mathbf{V}'\mathbf{b}$  (principal component regression coefficients)
    - $\mathbf{T} = \mathbf{U}\mathbf{S}$  (=  $\mathbf{M}\mathbf{V}$ ) (score matrix)
- Dimension reduction, include the first  $q$  principal components
  - $\mathbf{M} \approx \mathbf{U}_q\mathbf{S}_q\mathbf{V}_q'$
  - $\mathbf{T} = \mathbf{U}_q\mathbf{S}_q$  (=  $\mathbf{M}\mathbf{V}_q$ )
- Performing SVD is demanding for large datasets

# Chromosome-wise SVD on a core sample



# Chromosome-wise SVD on a core sample



# Single-step marker effects model (ssMEM)

- Fernando et al. (GSE 2016, 48:96)
- Compute expected genotypes for non-genotyped animals by solving:
  - $A^{22}\hat{M}_2 = -A^{21}M_1$
  - Total genotype matrix (genotyped and ungenotyped) is:
    - $M = \begin{bmatrix} M_1 \\ \hat{M}_2 \end{bmatrix}$
- ssMEM:
  - $y = ZMb + Z_2\epsilon + e$ 
    - where  $\epsilon \sim N(0, (A^{22})^{-1}\sigma_a^2)$
- ssMEM equations:
  - $$\begin{bmatrix} M'Z'ZM + I\rho\lambda & M'Z'Z_2 \\ Z_2'ZM & Z_2'Z_2 + A^{22}\lambda \end{bmatrix} \begin{bmatrix} \hat{b} \\ \hat{\epsilon} \end{bmatrix} = \begin{bmatrix} M'Z'y \\ Z_2'y \end{bmatrix}$$
  - where  $\lambda = \frac{\sigma_e^2}{\sigma_a^2}$

# Single-step principal component ridge-regression (ssPCRR)

- Compute expected scores for all non genotyped animals by solving:
  - $A^{22}\hat{C}_2 = -A^{21}C_1$  ( $C_1$  = approx. scores of genotyped)
  - Total score matrix (genotyped and ungenotyped) is now:  $C = \begin{bmatrix} C_1 \\ \hat{C}_2 \end{bmatrix}$
- ssPCRR model:
  - $y = ZCs + Z_2\epsilon + e$
- ssPCRR equations:
  - $$\begin{bmatrix} C'Z'ZC + I\rho\lambda & C'Z'Z_2 \\ Z_2'ZC & Z_2'Z_2 + A^{22}\lambda \end{bmatrix} \begin{bmatrix} \hat{s} \\ \hat{\epsilon} \end{bmatrix} = \begin{bmatrix} C'Z'y \\ Z_2'y \end{bmatrix}$$
- Genotyped EBV:
  - $\hat{a}_1 = C_1\hat{s}$
- Ungenotyped EBV
  - $\hat{a}_2 = C_2\hat{s} + \hat{\epsilon}$



# Simulation study

- Simulated population using QMSim (Sargolzaei and Schenkel, 2009)
  - 30 chromosomes of 100 cM
    - 24,259 SNP marker loci
    - 829 QTL
  - $h^2 = 0.25$
  - $N_e = 500$
  - 20,000 genotyped
  - 100,000 ungenotyped
  - All animals had own phenotype
- Chromosome-wise SVD
  - 2000 core animals
  - Number of chosen components set to explain >99% of genomic variation
- Block-iterative solver
- All analyses were run in a Julia environment (<https://julialang.org/>)

# Performance of models

- If full-scale SVD is performed
  - All models are equivalent and give identical results
- (Chromosome-wise) Reduced-dimension ssPCRR
  - EBV correlation to original ssGBLUP was  $>0.9999$
- Large-scale analysis
  - 4710 PC needed (157 per chromosome)
  - Setting up equation system ~ 4 minutes
  - Solving ~ 3 minutes
- Accuracies:
  - Genotyped: 0.90
  - Ungenotyped: 0.76

# Conclusions

- Large-scale genomic data from populations of limited  $N_e$ 
  - Few PC capture nearly all genetic variation
    - $\ll$  number of loci (dense data)
    - $\ll$  number of genotyped animals (large  $N$ )
- Fast SVD and dimension reduction
  - Smaller core sample
  - Parallel chromosome-wise SVD
- Single-step PC ridge regression (ssPCRR)
  - Very close approximation of the original ssGBLUP EBVs
  - Dimension of equation system greatly reduced
  - No need for inverse relationship matrices of genotyped animals

# Acknowledgements



Project no. 255297: "From whole genome sequence to precision breeding"

