# Progress in GWAS for large datasets with GBLUP and single-step GBLUP

Ignacy Misztal, Daniela Lourenco and Matias Bermann

# Specificity of plant and animal breeding

- Plants
  - Find genes in wild species
  - Introgress into inbred lines
  - Genetic evaluation of inbred crosses across environments
    - All crosses genotyped


- Animals
  - Selection usually within breeds and lines
  - Commercial animals purebreds or crossbreds
  - Many animals ungenotyped
  - Single-step GBLUP dominant methodology

# Single-step GBLUP –pedigree and genomic relationships combined

Matrix H (Legarra ,2009)

$$H = A + \begin{bmatrix} A_{12}A_{22}^{-1} & 0 \\ 0 & I \end{bmatrix} \begin{bmatrix} I \\ I \end{bmatrix} [G - A_{22}] [I \quad I] \begin{bmatrix} A_{22}^{-1}A_{21} & 0 \\ 0 & I \end{bmatrix}$$

Inverse of H (Aguilar et al., 2010)

G –genomic relationship matrix
1 –ungenotyped animals
2-genotyped animals

$$H^{-1} = A^{-1} + \begin{bmatrix} 0 & 0 \\ 0 & G^{-1} - A_{22}^{-1} \end{bmatrix}$$

Christensen and Lund, 2010
Boemcke et al., 2011

# ssGBLUP for Genome Wide Association Studies

- Large research interest in GWAS

- Limitations for current methods
  - Simple models
  - Single trait
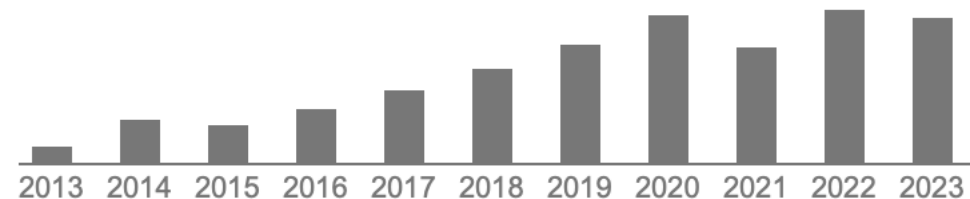  - Complicated if not all animals genotyped

Can ssGBLUP be used for GWAS?

Genome-wide association mapping including phenotypes
from relatives without genotypes

H. WANG[1]*, I. MISZTAL[1], I. AGUILAR[2], A. LEGARRA[3] AND W. M. MUIR[4]
[1] *Department of Animal and Dairy Science, University of Georgia, Athens, GA 30602-2771, USA*
[2] *Instituto Nacional de Investigación Agropecuaria, INIA Las Brujas, 90200 Canelones, Uruguay*
[3] *INRA, UR631 Station d'Amélioration Génétique des Animaux (SAGA), BP 52627, 32326 Castanet-Tolosan, France*
[4] *Department of Animal Science, Purdue University, West Lafayette, IN 47907-1151, USA*

Cited by 537

| 2013 | 2014 | 2015 | 2016 | 2017 | 2018 | 2019 | 2020 | 2021 | 2022 | 2023 |

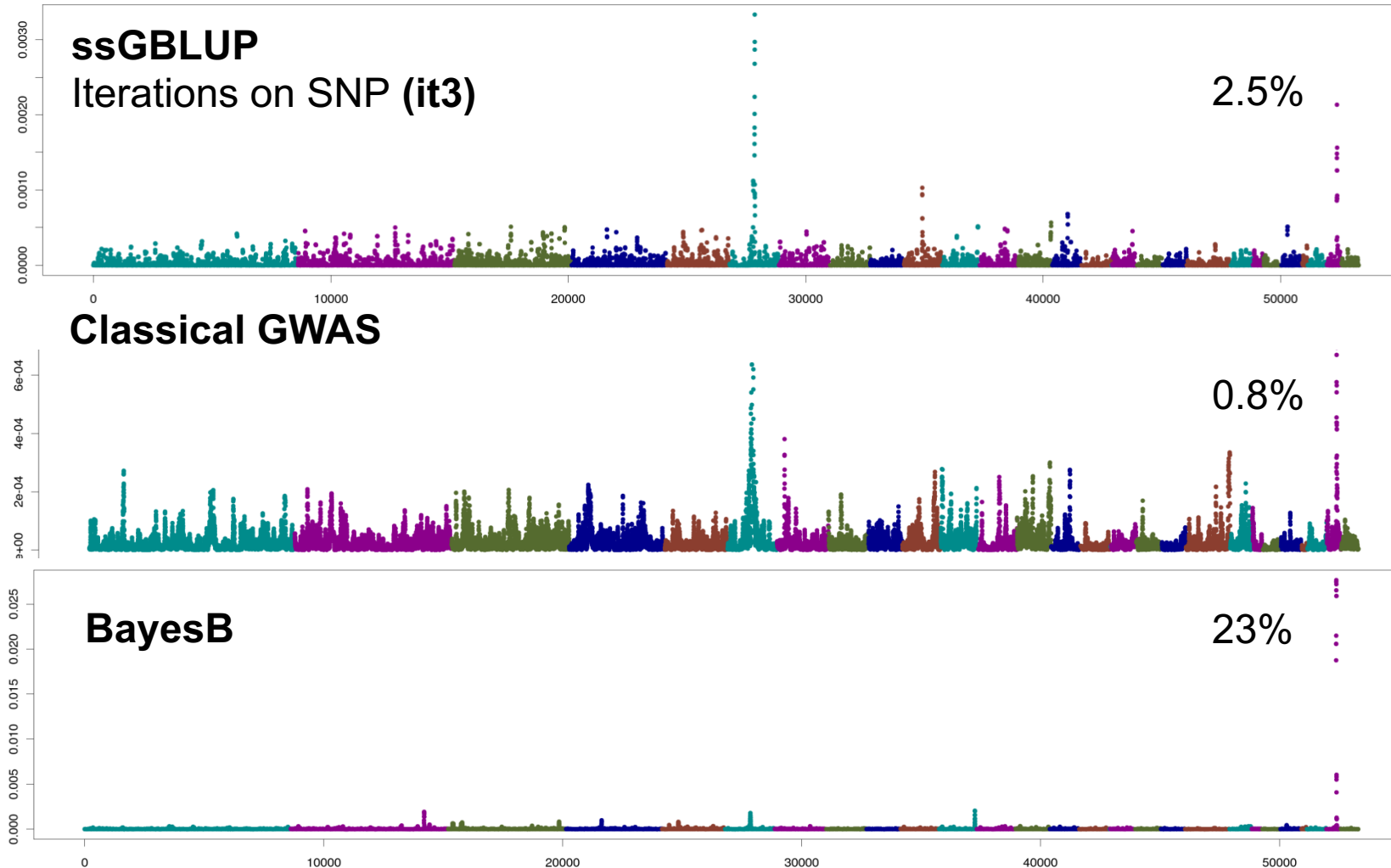# GWAS with ssGBLUP (Wang et al., 2012)

- Convert GEBV to SNP effects
- Estimate individual SNP variances
- Incorporate variances in G
- Possibly recompute GEBV and iterate

1. $D=I$
2. $G=ZDZ'/q$
3. Compute $a$
4. $u=DZ'/q \, G^{-1} a$
5. $d_i=2p_i(1-p_i)u_i^2$
6. $D=n \, D/tr(D)$
7. Loop to 2

Output as % of variance explained in a window

# Discrepancies in GWAS methods
## Chicken weight

# P-values for GWAS in (ss)GBLUP

$$pval_i = 2\left(1 - \Phi\left(\left|\frac{\widehat{snp}_i}{sd(\widehat{snp}_i)}\right|\right)\right) \text{ (Chen et al., 2017)}$$

If $sd(\widehat{snp}_i)$ approximately constant, Manhattan plots based on $|\widehat{snp}_i|$ and $pval_i$ similar

# Large data – APY algorithm

- Due to LD, genomic information compresses well: about 15k for cattle and about 5k for pigs and chicken

- APY algorithm: $u_{noncore} = P\ u_{core}, + \varepsilon$

- Number of core animals ~ equal to dimensionality

**Using recursion to compute the inverse of the genomic relationship matrix**

I. Misztal,*[1] A. Legarra,† and I. Aguilar‡
*Department of Animal and Dairy Science, University of Georgia, Athens 30602-2771
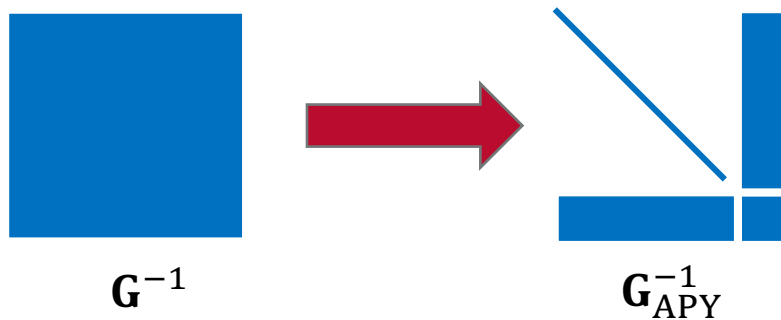†INRA, UR631-SAGA, BP 52627, 31326 Castanet-Tolosan Cedex, France
‡Instituto Nacional de Investigación Agropecuaria, Las Brujas 90200, Uruguay

**Inexpensive Computation of the Inverse of the Genomic Relationship Matrix in Populations with Small Effective Population Size**

Ignacy Misztal[1]
Animal and Dairy Science, University of Georgia, Athens, Georgia 30602
ORCID ID: 0000-0002-0382-1897 (I.M.)

$$\mathbf{A}^{-1} + \begin{pmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{G}^{-1} - \mathbf{A}_{22}^{-1} \end{pmatrix} \quad \Rightarrow \quad \mathbf{A}^{-1} + \begin{pmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{G}_{APY}^{-1} - \mathbf{A}_{22}^{-1} \end{pmatrix}$$

$\mathbf{G}^{-1}$

$\mathbf{G}_{APY}^{-1}$

# APY Single-step GWAS

- **Model**

$$y = W\alpha + Zu + \eta$$

- **Procedure**
    1. Calculate $Var(\boldsymbol{u})^{-1} = \boldsymbol{H}_{APY}^{-1}$
    2. Estimate variance components
    3. Calculate $\widehat{\boldsymbol{u}}_{2_c}$ and approximate $Var(\widehat{\boldsymbol{u}}_{2_c}) = \boldsymbol{G}_{cc} - \boldsymbol{C}^{\boldsymbol{u}_{2c}\boldsymbol{u}_{2c}}$
    4. For each marker:
        1. Calculate $\hat{b}_i = \boldsymbol{x}'_{c_i} \boldsymbol{G}_{cc}^{-1} \widehat{\boldsymbol{u}}_{\boldsymbol{2}}$
        2. Calculate $sd(\hat{b}_i) = \sqrt{\boldsymbol{x}'_{c_i} \boldsymbol{G}_{cc}^{-1}(\boldsymbol{G}_{cc} - \boldsymbol{C}^{\boldsymbol{u}_{2c}\boldsymbol{u}_{2c}})\boldsymbol{G}_{cc}^{-1}\boldsymbol{x}_{c_i}}$
        3. Calculate p-value as $pvalue_i = 1 - \Phi\left(\frac{\hat{b}_i}{sd(\hat{b}_i)}\right)$

Matias Bermann[1*], Daniela Lourenco[1], Natalia S. Forneris[2,3], Andres Legarra[4] and Ignacy Misztal[1]

**Matias Bermann,[1] Daniela Lourenco, and Ignacy Misztal**
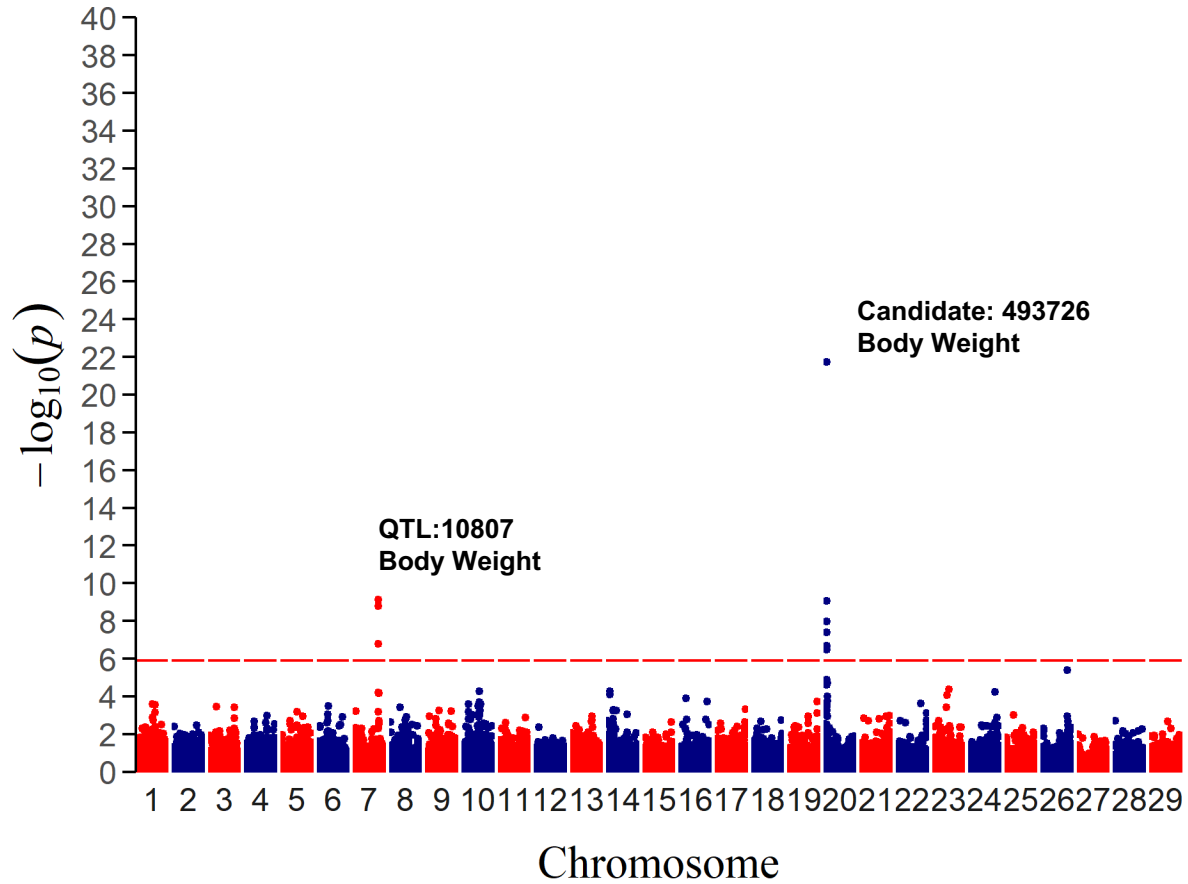
# Application example

- Post-weaning gain in American Angus

- 845,000 phenotypes

- 450,000 genotypes

- 1,570,000 animals in the pedigree

- ssGWAS (50k genotyped animals) vs. APY-ssGWAS (450k genotyped animals)

- We expect:
  - Higher power
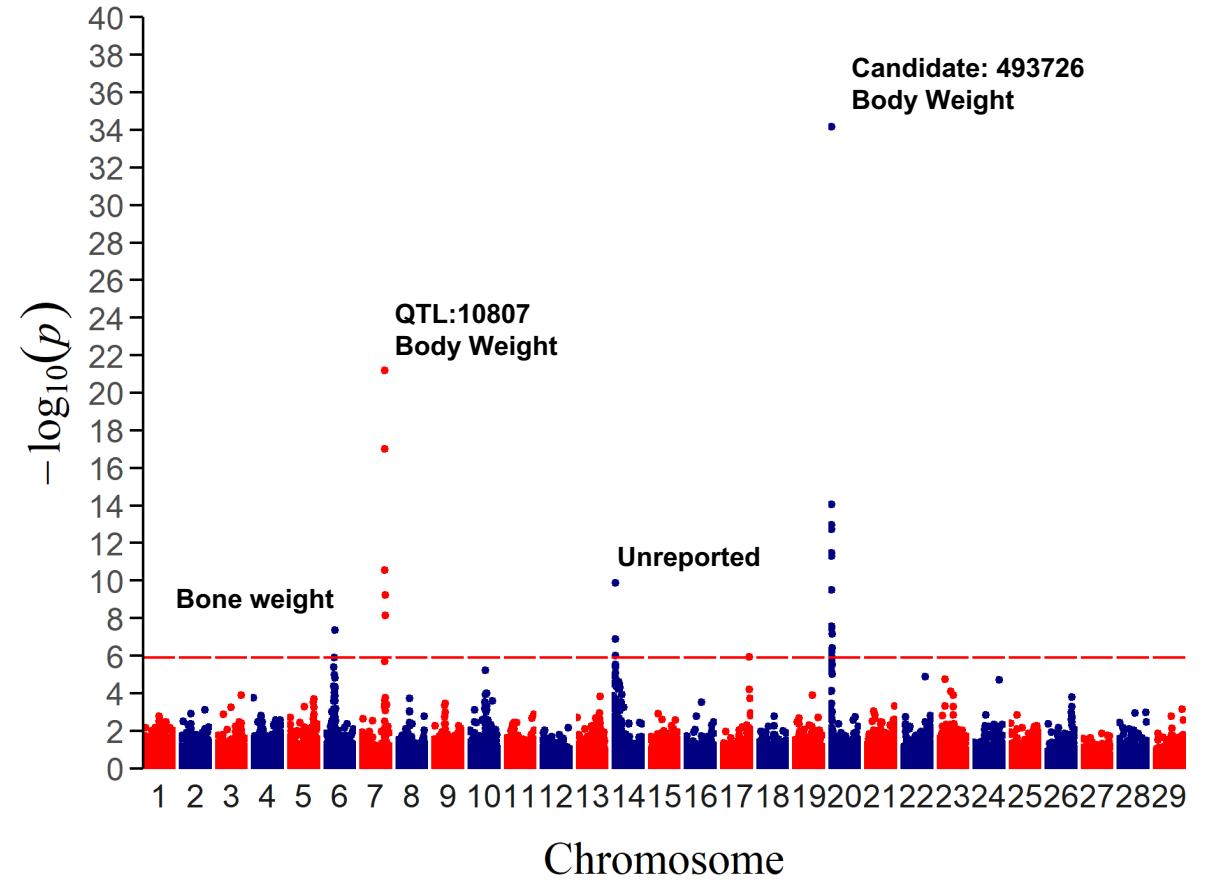  - Less noise
  - Less false-positives

Leite et al.
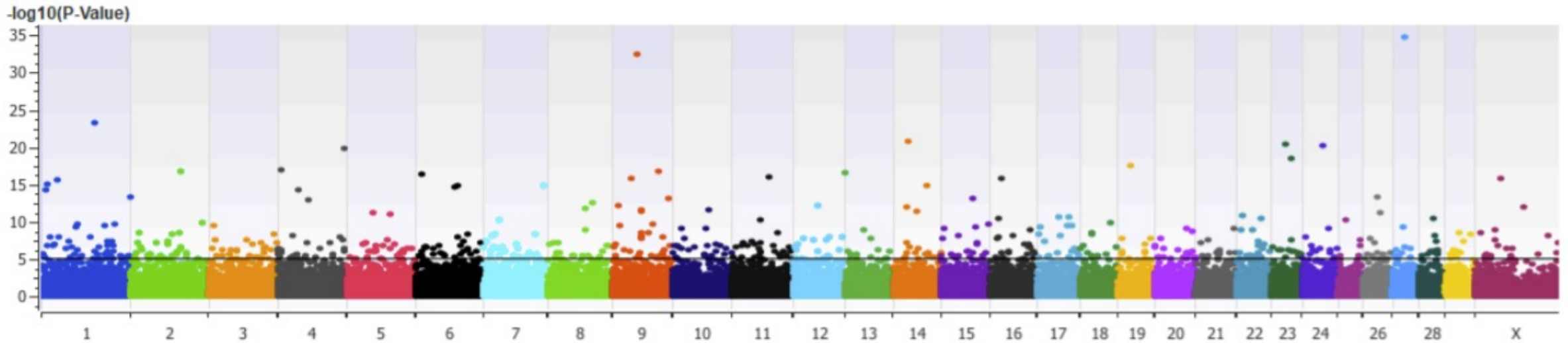(in progress)

**50k** genotyped animals

**500k** genotyped animals

# Questions with GWAS and predictions

- GWAS by
  - % of variance explained usually per 1Mb
  - p-values
- Few regions explain > 1% additive variance
- Lots of QTLs detected with small data sets
- Fewer QTLs detected with large data

# First conception rate on 2k Holstein heifers



Estimated heritability 36% (normally 1%)

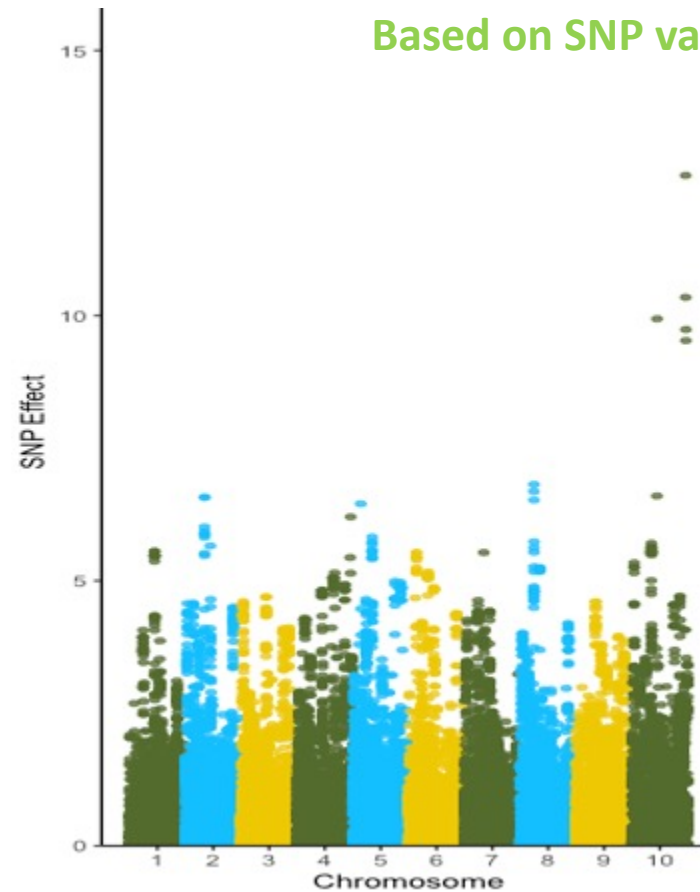Identified 146 unique loci at $p < 5 \times 10^{-8}$ level

Galliou et al., 2020, https://doi.org/10.3390/genes11070767

# Manhattan plots for simulated population with 100 identical equidistant QTNs



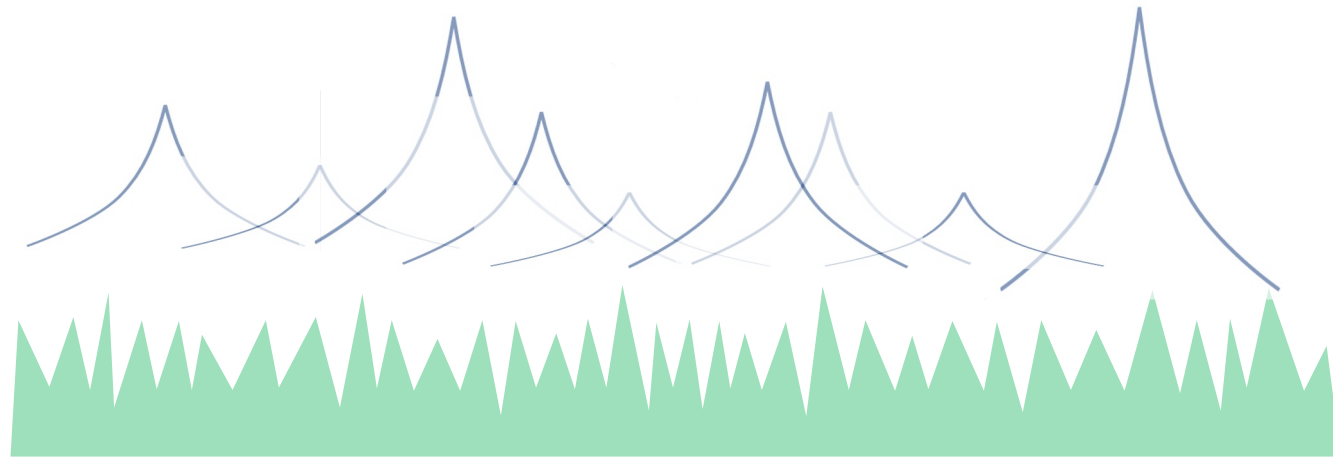Expectation

Based on SNP values

Based on p-values

Work started by Pocrnic et al. (2018)

# Plots averaged for 100 QTN



Pairwise linkage disequilibrium curve

~ 2 Mb for cattle
~ 5 Mb for pigs/chickens

~ 15 kb for humans

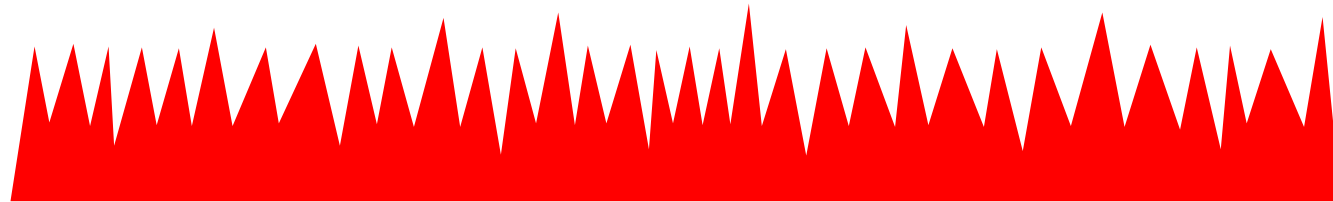1/Ne Morgans for 80% QTN variance
Ne - effective population size
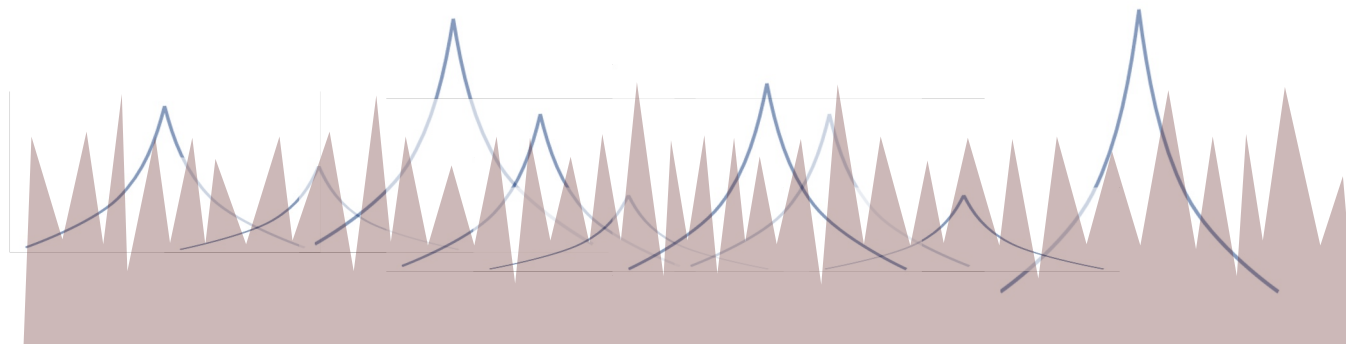
# What is Manhattan plot composed of?



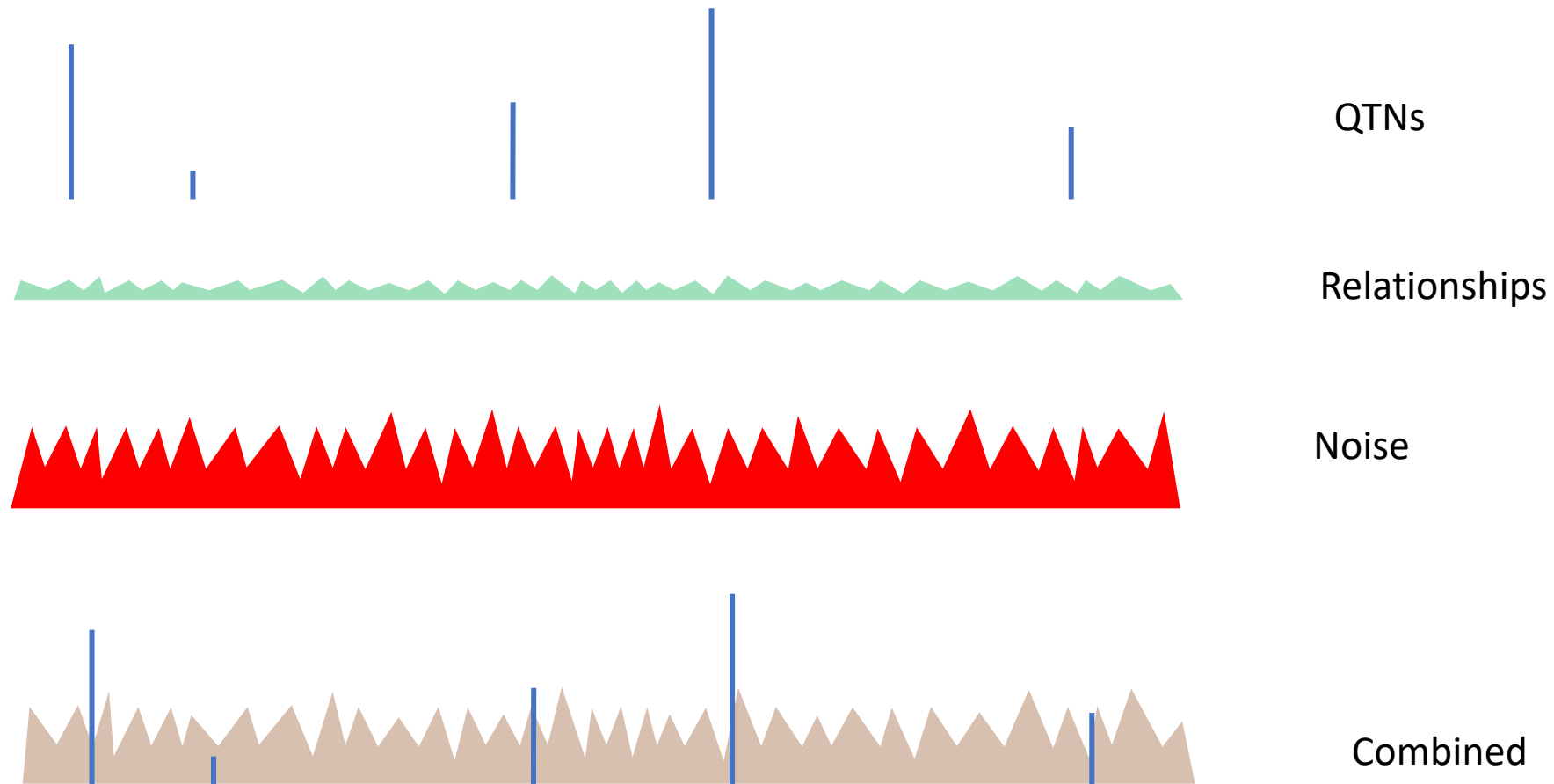QTNs — **Bigger with larger QTN and larger data**

Relationships

Noise — **Smaller with more data**

Combined
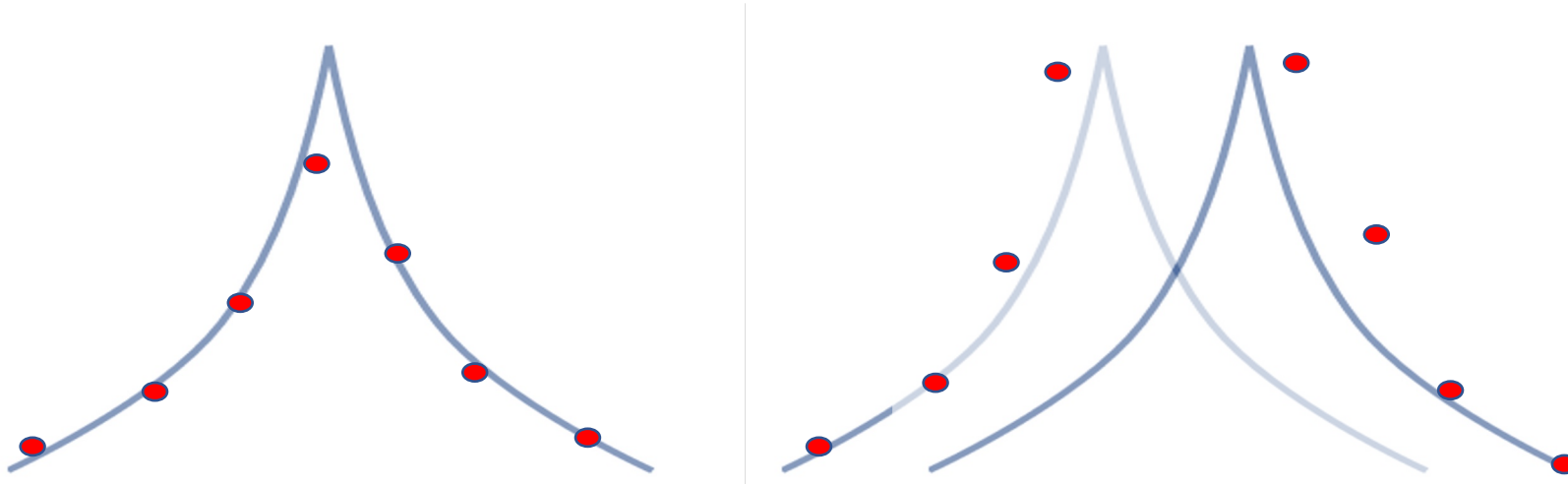
# Large effective population size



QTNs

Relationships

Noise

Combined

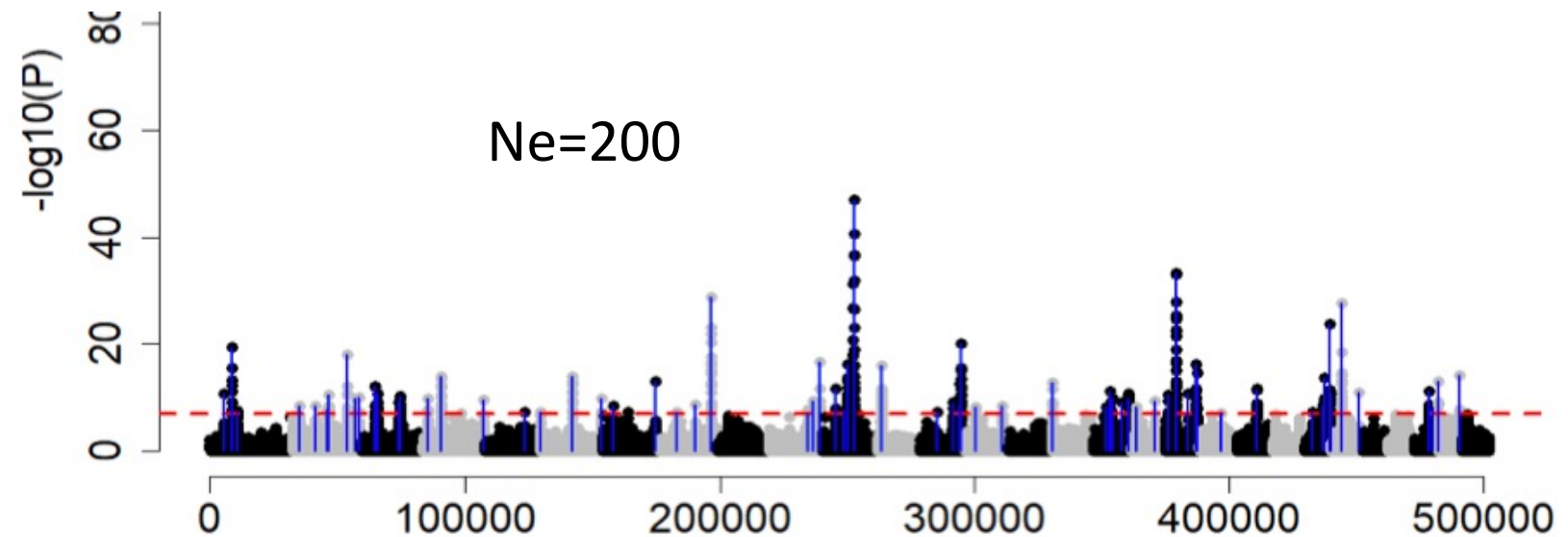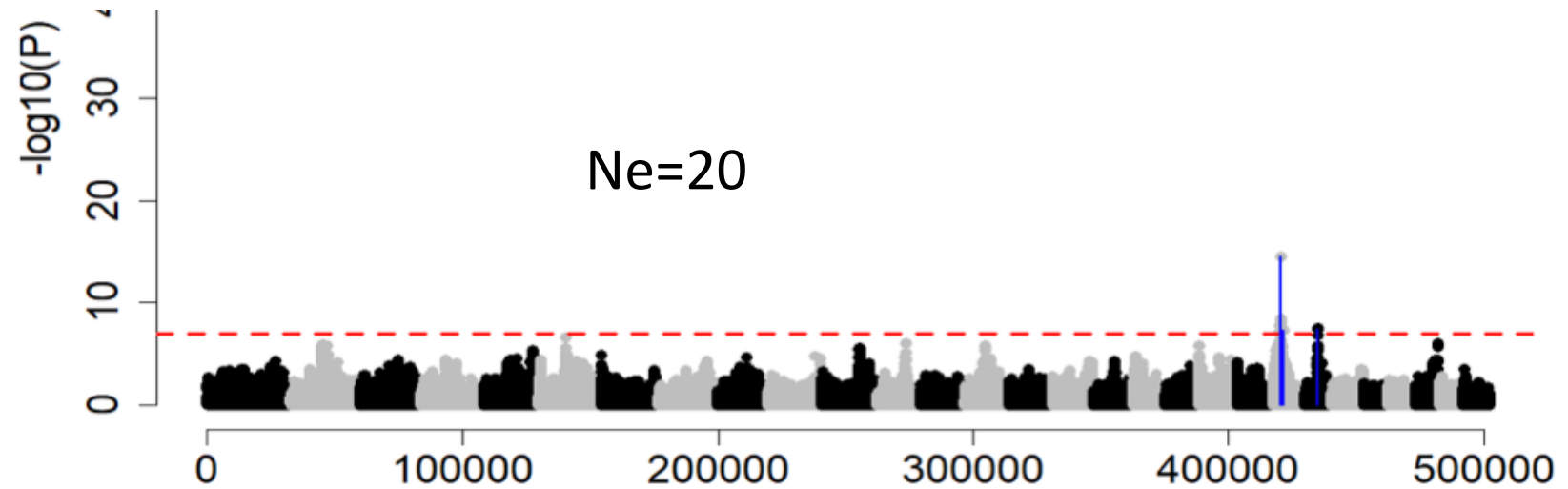# Why GBLUP accounts for QTN?



If 4 SNP per segment, 32 SNP  account for 80% of QTN variance

Need chip with 16 NeL SNP to mostly account for QTN

About 20k for pigs/broilers, 60k for cattle,  5m for humans
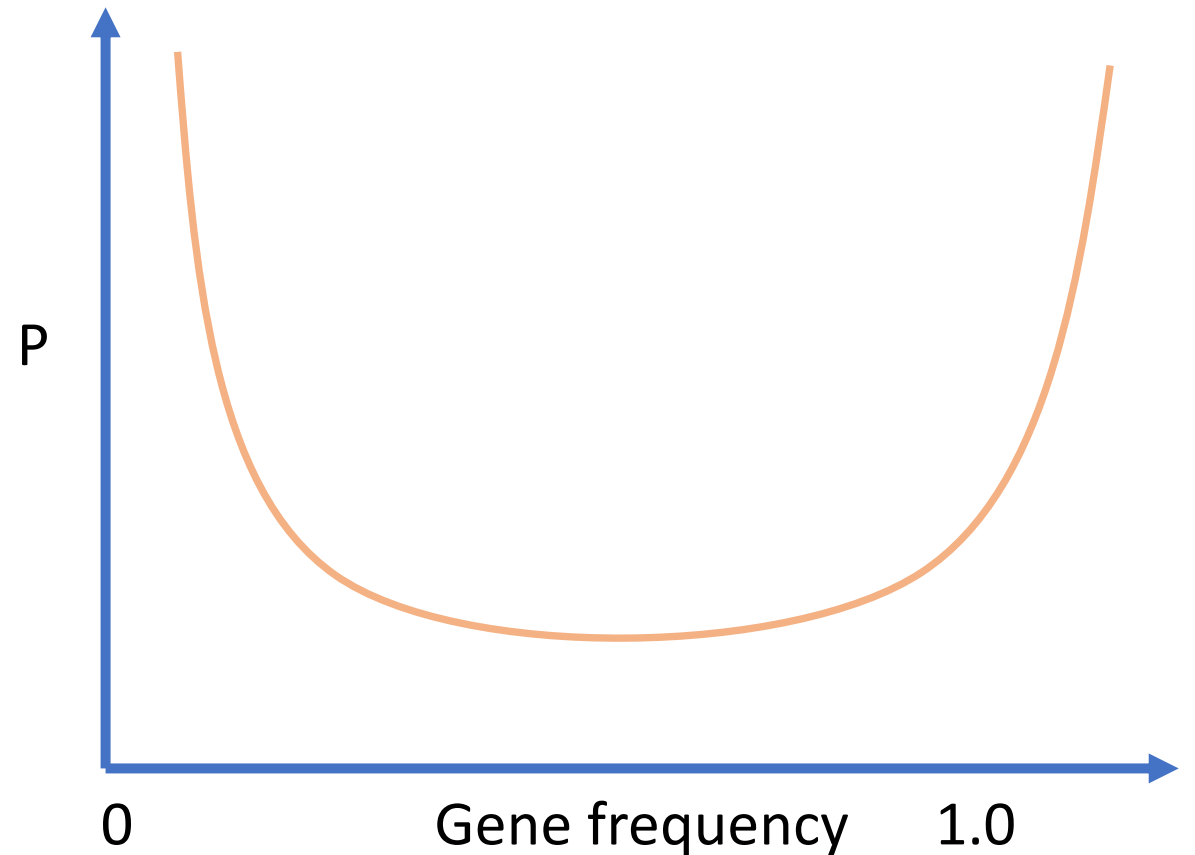
# Effective population size affects GWAS



Ne=20

Ne=200

Sungbong et al., 2021

# Why few QTN detected?

**AlphaSimR: an R package for breeding program simulations** 🔓
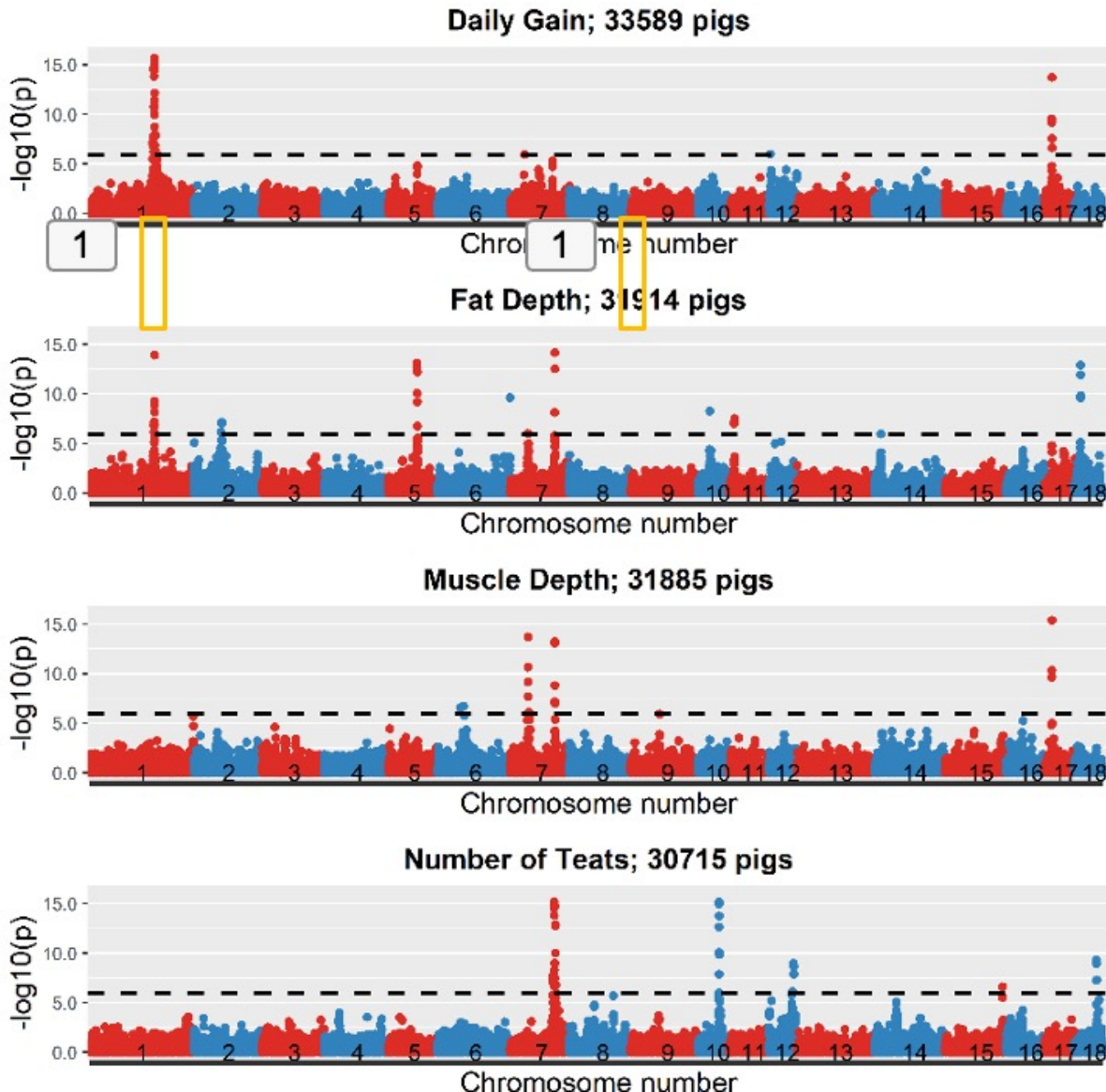
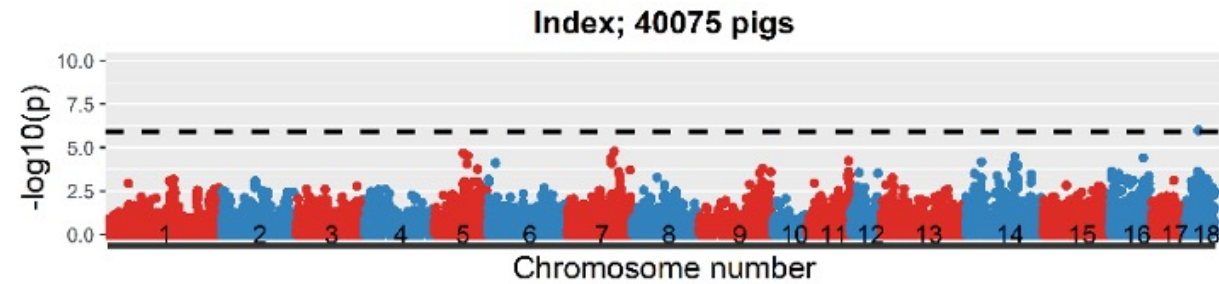R Chris Gaynor ✉, Gregor Gorjanc, John M Hickey

Only 20-30% QTN with p > 0.3

P

0      Gene frequency      1.0

# GWAS for various traits and index in pigs

Bijma, EAAP 23


Daily Gain; 33589 pigs


Fat Depth; 31914 pigs


Muscle Depth; 31885 pigs


Number of Teats; 30715 pigs

**Index**


Index; 40075 pigs

- Different peaks in different lines
- Antagonistic pleiotropy

# Conclusions

- GWAS in farm animals affected by small effective population size

- Optimal window size 1-2 Mb for Ne=100

- Large signals in GWAS due to QTN, relationships and noise (incl. Imputation)

- Large QTL show pleiotropy – QTL not visible in index

- GWAS by single-step GBLUP for any data size with option for p-values