

Single-step GWAS with p-values for large genotyped populations

Natalia Leite, Matias Bermann, Shogo Tsuruta,
Ignacy Misztal, **Daniela Lourenco**

June 28, 2023



UNIVERSITY OF
GEORGIA

College of Agricultural &
Environmental Sciences

ADSA Annual Meeting

June 25-28, 2023
Ottawa, Ontario, Canada

Equivalence ssGBLUP – ssSNPBLUP

$$\begin{bmatrix} \mathbf{X}'\mathbf{X} & \mathbf{X}'\mathbf{W} \\ \mathbf{W}'\mathbf{X} & \mathbf{W}'\mathbf{W} + \mathbf{H}^{-1} \frac{\sigma_e^2}{\sigma_u^2} \end{bmatrix} \begin{bmatrix} \hat{\boldsymbol{\beta}} \\ \hat{\mathbf{u}} \end{bmatrix} = \begin{bmatrix} \mathbf{X}'\mathbf{y} \\ \mathbf{W}'\mathbf{y} \end{bmatrix}$$

Aguilar et al. (2010)

$$\begin{bmatrix} \mathbf{X}'\mathbf{X} & \mathbf{X}'_g \mathbf{Z}_g \mathbf{M}_g & \mathbf{X}'_n \mathbf{Z}_n \\ \mathbf{M}'_g \mathbf{Z}'_g \mathbf{X}_g & \mathbf{Q} & \mathbf{M}'_g \mathbf{A}^{gn} \frac{\sigma_e^2}{\sigma_g^2} \\ \mathbf{Z}'_n \mathbf{X}_n & \mathbf{A}^{ng} \mathbf{M}_g \frac{\sigma_e^2}{\sigma_g^2} & \mathbf{Z}'_n \mathbf{Z}_n + \mathbf{A}^{nn} \frac{\sigma_e^2}{\sigma_g^2} \end{bmatrix} \begin{bmatrix} \hat{\boldsymbol{\beta}} \\ \hat{\boldsymbol{\alpha}} \\ \hat{\mathbf{u}}_n \end{bmatrix} = \begin{bmatrix} \mathbf{X}'\mathbf{y} \\ \mathbf{M}'_g \mathbf{Z}'_g \mathbf{y}_g \\ \mathbf{Z}'_n \mathbf{y}_n \end{bmatrix}$$

Fernando et al. (2016)

- Equivalent models under same assumptions and data

- $\hat{\mathbf{u}} = \mathbf{Z}\hat{\mathbf{a}}$

- $\hat{\mathbf{a}}|\hat{\mathbf{u}} = k\mathbf{Z}'\mathbf{G}^{-1}\hat{\mathbf{u}}$

- $\text{Var}(\hat{\mathbf{a}}|\hat{\mathbf{u}}) = k\mathbf{Z}'\mathbf{G}^{-1}(\mathbf{G} - \mathbf{C}^{\mathbf{u}_2\mathbf{u}_2})\mathbf{G}^{-1}\mathbf{Z}k$

VanRaden (2008)

Stranden and Garrick (2009)

Guladron-Duarte et al. (2014)

- $p\text{-value}_i = 2 \left(1 - \Phi \left(\left| \frac{\hat{a}_i}{sd(\hat{a}_i)} \right| \right) \right) \rightarrow \text{ssGWAS}$

Aguilar et al. Genet Sel Evol (2019) 51:28
<https://doi.org/10.1186/s12711-019-0469-3>

GSE Genetics
 Selection
 Evolution

SHORT COMMUNICATION

Open Access

Frequentist p-values for large-scale single step genome-wide association, with an application to birth weight in American Angus cattle

Ignacio Aguilar¹, Andres Legarra^{2*}, Fernando Cardoso^{3,4}, Yutaka Masuda⁵, Daniela Lourenco⁵ and Ignacy Misztal⁶

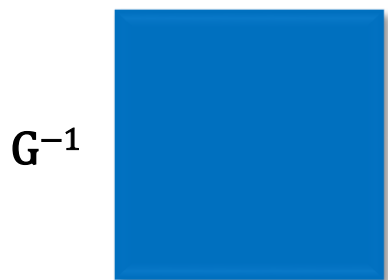
Equivalence enables ssGWAS

- Why ssGWAS?
- Assumption of single-marker GWAS: Genotyped individuals have phenotypes
- Animal populations: genotypes and phenotypes may not be on the same individuals
 - Deregressed EBV: are biased
- ssGWAS
 - All data on genotyped and non-genotyped individuals
 - Multi-trait models to accommodate correlations
- Negative aspect of ssGWAS
 - Heavy computations for p-values -> same limitation as REML

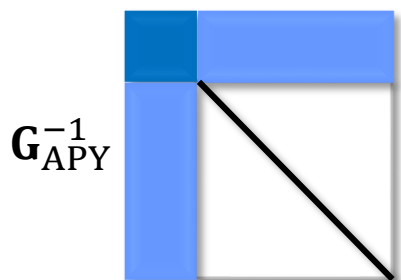
Algorithm for Proven and Young (APY)

- Realized relationship matrix in ssGBLUP

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



Dense $\rightarrow u_i | u_1 + u_2 + u_3, \dots, u_{i-1} = \sum_{j=1}^{n-1} p_{ij} u_j + \varepsilon_i$



Sparse $\rightarrow u_i | u_{c1} + u_{c2} + u_{c3}, \dots, u_{ci} = \sum_{j=1}^c p_{ij} u_j + \varepsilon_i$

Condition on a set of features or animals = CORE animals

Misztal et al. (2014)
Fragomeni et al. (2015)
Lourenco et al. (2015)



- \mathbf{G}_{APY}^{-1} sparse
- Efficient computations

Masuda et al. (2016)

Equivalence APY ssGBLUP – ssSNPBLUP



JOURNAL ARTICLE

Indirect predictions with a large number of genotyped animals using the algorithm for proven and young

Andre L.S. Garcia, Yutaka Masuda, Shogo Tsuruta, Stephen Miller, Ignacy Misztal, Daniela Lourenco

Journal of Animal Science, Volume 98, Issue 6, June 2020, skaa154, <https://doi.org/10.1093/jas/skaa154>

Garcia et al. *Genetics Selection Evolution* (2022) 54:66 <https://doi.org/10.1186/s12711-022-00752-4>



RESEARCH ARTICLE

Open Access

Theoretical accuracy for indirect predictions based on SNP effects from single-step GB LUP

Andre Garcia, Ignacio Aguilar, Andres Legarra, Shogo Tsuruta, Ignacy Misztal and Daniela Lourenco

- If using APY in ssGBLUP
 - Numerical equivalence
 - $\hat{\mathbf{u}} = \mathbf{Z}\hat{\mathbf{a}}$
 - $\hat{\mathbf{a}}|\hat{\mathbf{u}} = k\mathbf{Z}'\mathbf{G}_{\text{APY}}^{-1}\hat{\mathbf{u}}$
 - $\text{Var}(\hat{\mathbf{a}}|\hat{\mathbf{u}}) = k\mathbf{Z}'\mathbf{G}_{\text{APY}}^{-1}(\mathbf{G} - \mathbf{C}^{\mathbf{u}_2\mathbf{u}_2})\mathbf{G}_{\text{APY}}^{-1}\mathbf{Z}k$

↳ Function of all genotyped animals



Bermann et al. *Genetics Selection Evolution* (2022) 54:52 <https://doi.org/10.1186/s12711-022-00741-7>



RESEARCH ARTICLE

Open Access

On the equivalence between marker effect models and breeding value models and direct genomic values with the Algorithm for Proven and Young

Matias Bermann, Daniela Lourenco, Natalia S. Forneris, Andres Legarra and Ignacy Misztal

- If using APY in ssGBLUP
 - Equivalent APY ssSNPBLUP model
 - $\hat{\mathbf{u}} = \mathbf{Z}^+\hat{\mathbf{a}}$
 - $\hat{\mathbf{a}}|\hat{\mathbf{u}} = k\mathbf{Z}^+\mathbf{G}_{\text{APY}}^{-1}\hat{\mathbf{u}} = k\mathbf{Z}'_c\mathbf{G}_{\text{CC}}^{-1}\hat{\mathbf{u}}_c$
 - $\text{Var}(\hat{\mathbf{a}}|\hat{\mathbf{u}}) = k\mathbf{Z}'_c\mathbf{G}_{\text{CC}}^{-1}(\mathbf{G}_{\text{CC}} - \mathbf{C}^{\mathbf{u}_2\mathbf{u}_2})\mathbf{G}_{\text{CC}}^{-1}\mathbf{Z}_ck$

↳ Function of CORE animals

Equivalence APY ssGBLUP – ssSNPBLUP

- If using APY in ssGBLUP
 - Equivalent APY ssSNPBLUP model
 - $\hat{\mathbf{u}} = \mathbf{Z}^{\dagger} \hat{\mathbf{a}}$
 - $\hat{\mathbf{a}} | \hat{\mathbf{u}} = k \mathbf{Z}^{\dagger} \mathbf{G}_{\text{APY}}^{-1} \hat{\mathbf{u}} = k \mathbf{Z}'_c \mathbf{G}_{\text{CC}}^{-1} \hat{\mathbf{u}}_c$
 - $\text{Var}(\hat{\mathbf{a}} | \hat{\mathbf{u}}) = k \mathbf{Z}'_c \mathbf{G}_{\text{CC}}^{-1} (\mathbf{G}_{\text{CC}} - \mathbf{C}^{\mathbf{u}_{2c} \mathbf{u}_{2c}}) \mathbf{G}_{\text{CC}}^{-1} \mathbf{Z}_c k$

$\mathbf{C}^{\mathbf{u}_{2c} \mathbf{u}_{2c}}$

?

- Exact - Inverse of the LHS of MME
- Approximating reliabilities of GEBV

On the equivalence between marker effect models and breeding value models and direct genomic values with the Algorithm for Proven and Young




Matias Bermann^{1*}, Daniela Lourenco¹, Natalia S. Forneris^{2,3}, Andres Legarra⁴ and Ignacy Misztal¹

GEBV are published with reliability

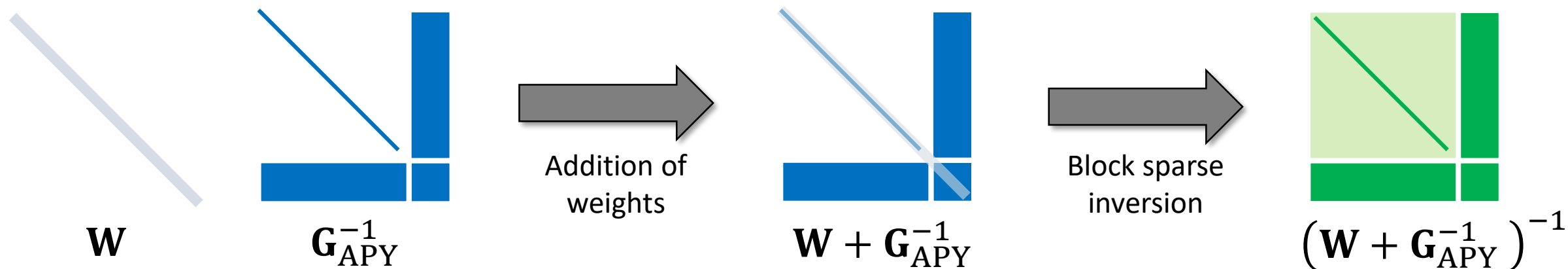
- Reliability based on PEV
 - Approximated for large populations
 - Weights based on approximations
 - Block sparse inversion with APY

JOURNAL ARTICLE

Efficient approximation of reliabilities for single-step genomic best linear unbiased predictor models with the Algorithm for Proven and Young 

Matias Bermann , Daniela Lourenco, Ignacy Misztal

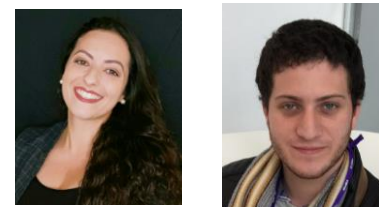
Journal of Animal Science, Volume 100, Issue 1, January 2022, skab353,
<https://doi.org/10.1093/jas/skab353>



$$diag(W + G_{APY}^{-1})^{-1} = \begin{matrix} diag((W_{nn} + M_{nn}^{-1})^{-1} + (W_{nn} + M_{nn}^{-1})^{-1} G^{nc} (W_{cc} + G^{cc} - G^{cn} (W_{nn} + M_{nn}^{-1})^{-1} G^{nc})^{-1} G^{cn} (W_{nn} + M_{nn}^{-1})^{-1}) \\ \\ diag((W_{cc} + G^{cc} - G^{cn} (W_{nn} + M_{nn}^{-1})^{-1} G^{nc})^{-1}) \end{matrix}$$

Single-step GWAS – Many genotypes

- Genomic evaluation process
 - GEBV using APY ssGBLUP + reliability using block sparse inversion



Leite, Bermann et al.
(in preparation)

- $\hat{\mathbf{a}}|\hat{\mathbf{u}} = \mathbf{k}\mathbf{Z}'_c\mathbf{G}_{cc}^{-1}\hat{\mathbf{u}}_c$
- $\mathbf{C}^{\mathbf{u}_{2c}\mathbf{u}_{2c}} = \left(\mathbf{W} + \frac{\sigma_e^2}{\sigma_u^2}\mathbf{G}_{APY}^{-1}\right)^{-1}$
- $\text{Var}(\hat{\mathbf{a}}|\hat{\mathbf{u}}) = \mathbf{k}\mathbf{Z}'_c\mathbf{G}_{cc}^{-1}(\mathbf{G}_{cc} - \mathbf{C}^{\mathbf{u}_{2c}\mathbf{u}_{2c}})\mathbf{G}_{cc}^{-1}\mathbf{Z}_c\mathbf{k}$

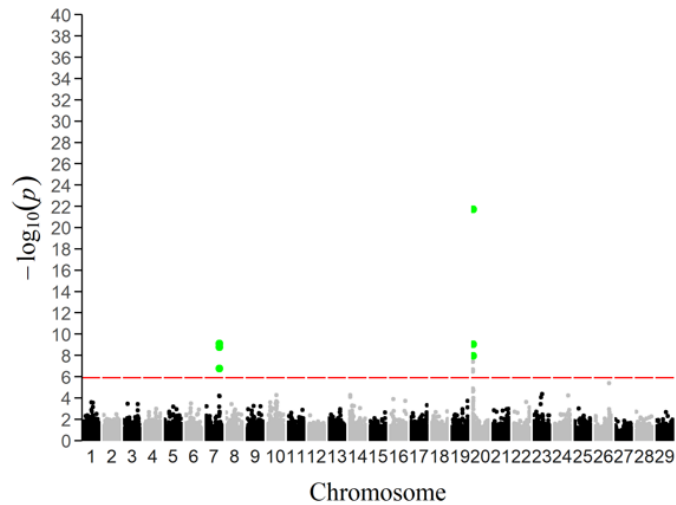
$$p\text{-value}_i = 2 \left(1 - \Phi \left(\left| \frac{\hat{a}_i}{sd(\hat{a}_i)} \right| \right) \right)$$

- Initial tests – AGI data
 - 845k phenotypes for post-weaning gain
 - 50k genotyped + 1.58M pedigree
 - 450k genotyped (13k core + 437k noncore) + 1.8M pedigree

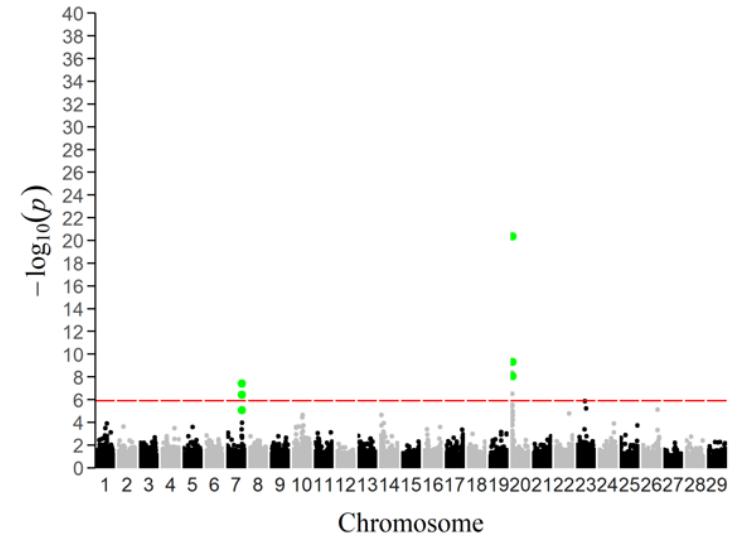


Single-step GWAS – Scenarios

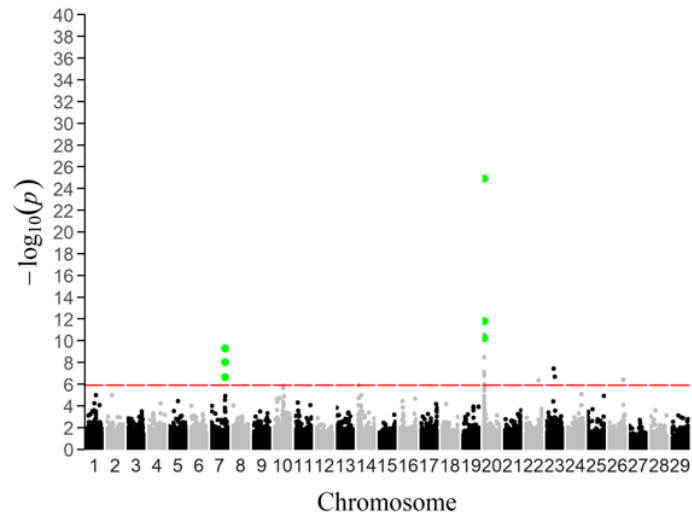
Exact - G^{-1} 50k



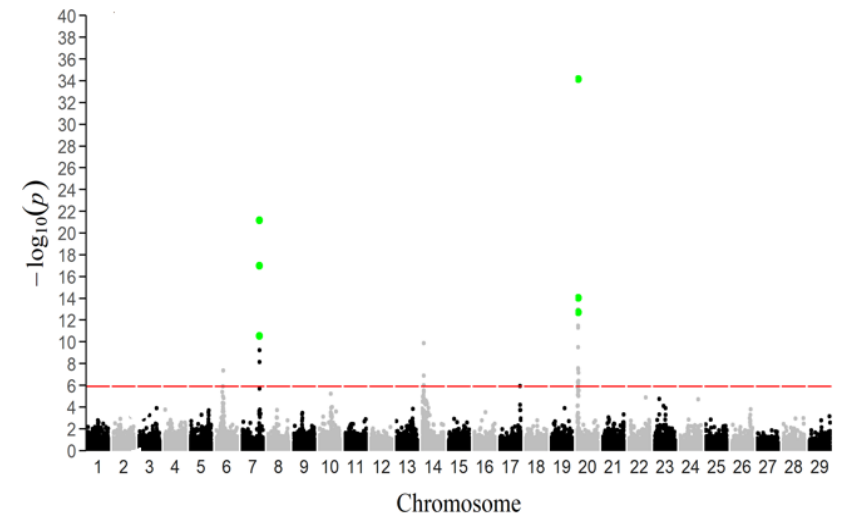
Exact - G_{APY}^{-1} 50k



Approximation - G_{APY}^{-1} 50k



Approximation - G_{APY}^{-1} 450k



Single-step GWAS – Computations

Method	Elapsed time, h:min*	Peak memory, GB*
Exact - G^{-1} 50k	106:46	159.66
Exact - G_{APY}^{-1} 50k	110:59	178.30
Approx - G_{APY}^{-1} 50k	2:50	16.62
Approx - G_{APY}^{-1} 450k	24:28	87.64

- Computing cost still high
- Eliminates the limitation on the amount of data for ssGWAS

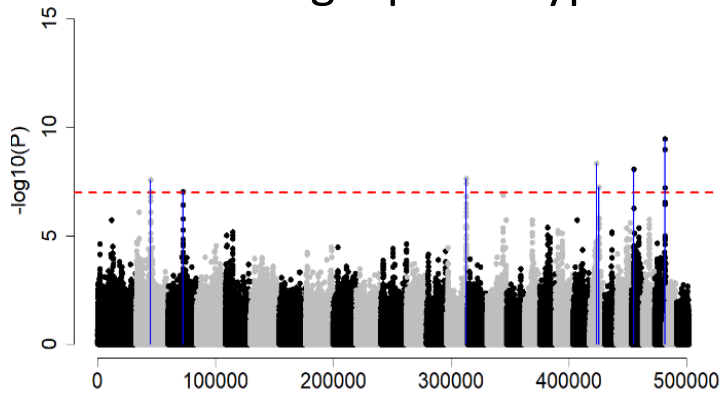
Adding more data in GWAS

- Adding more genotypes and phenotypes
 - Less noise
 - Much higher resolution

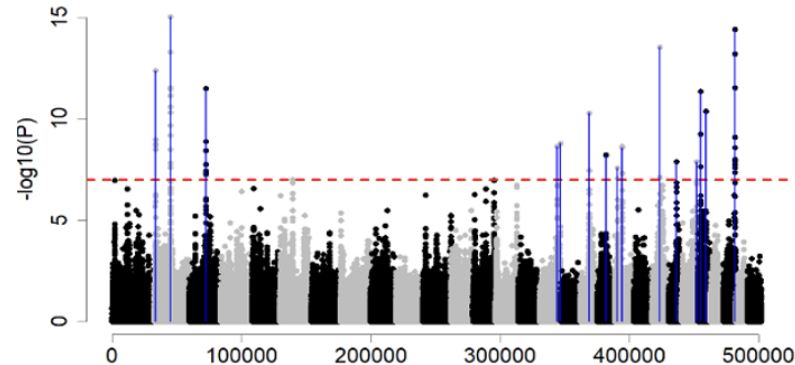


Jang et al.
(accepted)

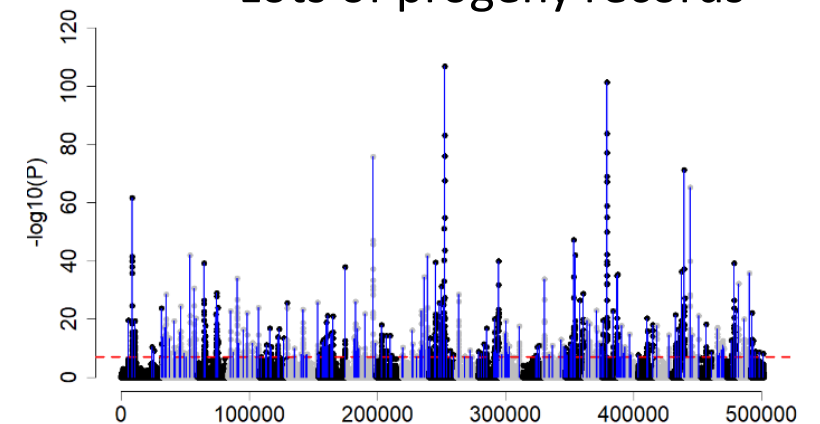
15,800 genotyped animals
+
Single phenotypes



30,000 genotyped animals
+
Single phenotypes



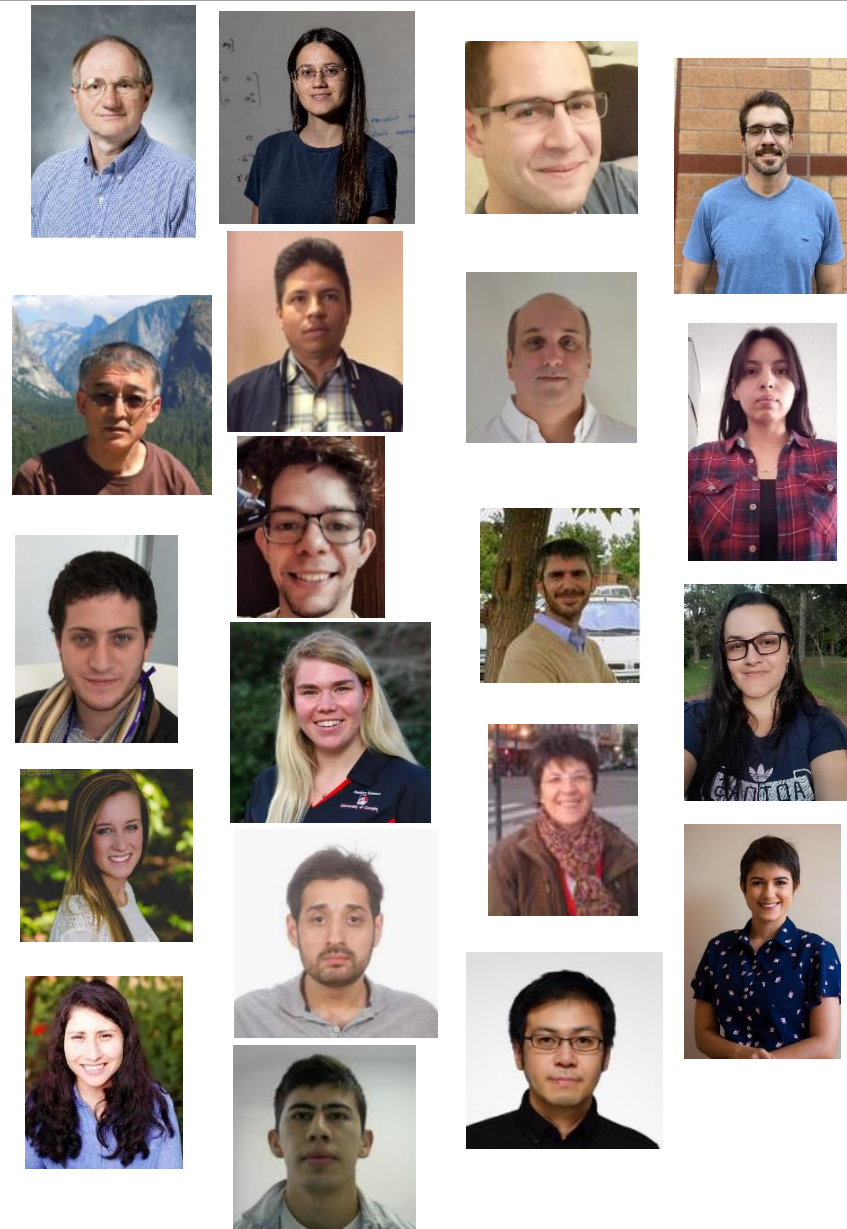
30,000 genotyped animals
+
Lots of progeny records



Take home messages

- ssGWAS allows for using all data
 - Phenotypes
 - Pedigree
 - Genotypes
- Virtually any number of genotyped animals
 - Improvement in computing time
- Possible because of the limited dimensionality of genomic information
 - Depends on the quality of approximated reliabilities
 - Already implemented in BLUPF90

UGA AB&G team



USDA United States Department of Agriculture
Agricultural Research Service

Warmwater Aquaculture Research Unit

USDA Agricultural Research Service
U.S. DEPARTMENT OF AGRICULTURE

Cool and Cold Water Aquaculture Research

