



# Efficient Estimation of Polygenic Effects via Multivariate Ridge Regression

**Alencar Xavier**

Research Scientist at Corteva Biostatistics  
Adjunct professor at Purdue University

**David Habier**

Sr. Research Scientist at Corteva Biostatistics

# Outline

## 1. Introduction

- Rationale and statistical model

## 2. Coefficients

- Univariate
- Multivariate

## 3. Variances

- Univariate
- Multivariate

## 4. Simulations

- Study 1: Comparison to REML in small balanced data
- Study 2: Performance in large unbalanced data
- Limitations and other considerations

## 5. Conclusion

# 1. Introduction

- Rationale and statistical model

## 2. Coefficients

- Univariate
- Multivariate

## 3. Variances

- Univariate
- Multivariate

## 4. Simulations

- Study 1: Comparison to REML in small balanced data
- Study 2: Performance in large unbalanced data
- Limitations and other considerations

## 5. Conclusion

# Rationale

- Single-trait models for genomic prediction in plant breeding are well-established (e.g. GBLUP and BayesB)
- Phenotypes come from multiple locations, years, and quantitative traits; and most traits have genetically correlated breeding values

# Rationale

- **Complex GxE / multi-trait patterns** (= higher accuracy)
- **Assess new phenomic traits** (e.g. canopy coverage in soy)
- **Computationally PROHIBITIVE\***

\* Zhou, X., & Stephens, M. (2014). Efficient multivariate linear mixed model algorithms for genome-wide association studies. *Nature methods*, 11(4), 407-409.

# Why would multivariate be any better?

Simple (bivariate) model:

INFORMATION GAIN

$$y = g + e$$

$$Var \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} \sigma_{a_1}^2 & \sigma_{a_{12}} \\ \sigma_{a_{12}} & \sigma_{a_2}^2 \end{bmatrix} + \begin{bmatrix} \sigma_{e_1}^2 & \sigma_{e_{12}} \\ \sigma_{e_{12}} & \sigma_{e_2}^2 \end{bmatrix}$$

# Why marker ridge regression?

1. Regression-type models are easy to store and use for prediction
2. Compatible with the multi-stage<sup>1,2</sup> framework
3. Well-known properties: Gaussian, additive, and equivalent to GBLUP
4. No need to build and invert G matrix (which is not always positive definite)
5. Provides covariance components for meaningful statistics:
  - Heritability, reliability, accuracy, genetic correlations, selection indexes, correlated response

1. Smith, A., Cullis, B., and Gilmour, A. (2001). Applications: the analysis of crop variety evaluation data in Australia. Australian & New Zealand Journal of Statistics, 43(2), 129-145.  
2. Mohring, J, and H-P Piepho, (2009) Comparison of weighting in two-stage analysis of plant breeding trials. Crop Sci. 49: 1977–1988.

# Statistical model

$$y = \mu + \mathbf{Z}\beta + e \quad (1)$$

- Where  $y = \{y_1, y_2, \dots, y_K\}$ ,  $\mu = \{\mu_1, \mu_2, \dots, \mu_K\}$ ,  $\beta = \{\beta_1, \beta_2, \dots, \beta_K\}$ ,  
 $e = \{e_1, e_2, \dots, e_K\}$ ,  $\mathbf{Z} = \text{BlockDiag}\{\mathbf{Z}_1, \mathbf{Z}_2, \dots, \mathbf{Z}_K\}$
- Variances:

$$\Sigma_{\beta} = \begin{bmatrix} \sigma_{\beta(1)}^2 & \dots & \sigma_{\beta(1,K)} \\ \vdots & \ddots & \vdots \\ \sigma_{\beta(K,1)} & \dots & \sigma_{\beta(K)}^2 \end{bmatrix} \quad \text{and} \quad \Sigma_e = \begin{bmatrix} \sigma_{e(1)}^2 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sigma_{e(K)}^2 \end{bmatrix}$$



# Corresponding mixed model equation

Under the traditional framework, the mixed-model equations required to solve the multivariate ridge regression (eq. 1) can be written as follows:

$$\begin{bmatrix} \mathbf{1}'_1 \mathbf{1}_1 \sigma_{e_1}^{-2} & \dots & 0 & \mathbf{1}'_1 \mathbf{Z}_1 \sigma_{e_1}^{-2} & \dots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & \mathbf{1}'_K \mathbf{1}_K \sigma_{e_K}^{-2} & 0 & \dots & \mathbf{1}'_K \mathbf{Z}_K \sigma_{e_K}^{-2} \\ \mathbf{Z}'_1 \mathbf{1}_1 \sigma_{e_1}^{-2} & \dots & 0 & \mathbf{Z}'_1 \mathbf{Z}_1 \sigma_{e_1}^{-2} + \mathbf{I}_m \sigma_{\beta}^{11} & \dots & \mathbf{I}_m \sigma_{\beta}^{1K} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & \mathbf{Z}'_K \mathbf{1}_K \sigma_{e_K}^{-2} & \mathbf{I}_m \sigma_{\beta}^{K1} & \vdots & \mathbf{Z}'_K \mathbf{Z}_K \sigma_{e_K}^{-2} + \mathbf{I}_m \sigma_{\beta}^{KK} \end{bmatrix} \begin{bmatrix} \hat{\mu}_1 \\ \vdots \\ \hat{\mu}_k \\ \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_K \end{bmatrix} = \begin{bmatrix} \sigma_{e_1}^{-2} \mathbf{1}'_1 \mathbf{y}_1 \\ \vdots \\ \sigma_{e_K}^{-2} \mathbf{1}'_K \mathbf{y}_K \\ \sigma_{e_1}^{-2} \mathbf{Z}'_1 \mathbf{y}_1 \\ \vdots \\ \sigma_{e_K}^{-2} \mathbf{Z}'_K \mathbf{y}_K \end{bmatrix} \quad (2)$$

where  $\sigma_{\beta}^{ij}$  is the element at position  $ij$  of  $\Sigma_{\beta}^{-1}$ . This setup involves storing  $K$  times the cross-product or marker scores ( $\mathbf{Z}'_k \mathbf{Z}_k$ ), each with dimension  $m \times m$ .

Moreover, this **huge** matrix must be **inverted** for the estimation of covariance components:  $\hat{\Sigma}_{\beta(i,j)} = m^{-1}[\hat{\beta}'_i \hat{\beta}_j + \text{tr}(\mathbf{C}^{ij})]$

Computing very large multivariate models is **impossible**

**unless...**



# 1. Introduction

- Rationale and statistical model

## 2. Coefficients

- Univariate
- Multivariate

## 3. Variances

- Univariate
- Multivariate

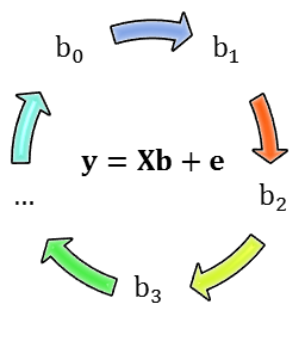
## 4. Simulations

- Study 1: Comparison to REML in small balanced data
- Study 2: Performance in large unbalanced data
- Limitations and other considerations

## 5. Conclusion

# Coefficients for univariate model

1. Whole-genome regression (e.g. BayesA) rely on the *Gauss-Seidel* method <sup>1</sup>
2. GS has only two steps, whereas coordinate descent has three <sup>2</sup>
3. It avoids building the systems of equations altogether!!
4. Estimates one marker effects, then uses residuals to update the next effect

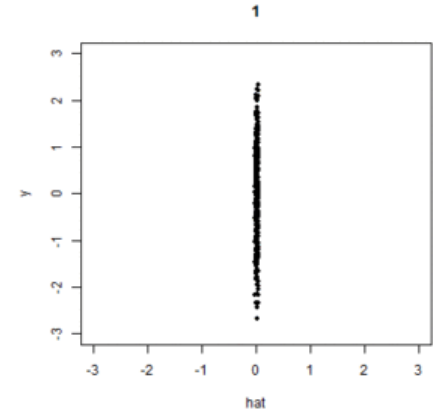
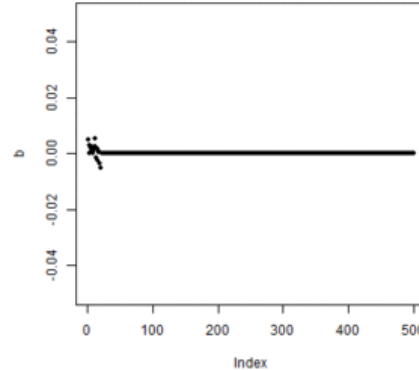


for j in 1:p {

$$\hat{b}_j^{t+1} = \frac{x_j' \hat{e}^t + x_j' x_j \hat{b}_j^t}{x_j' x_j + \lambda}$$

$$\hat{e}^{t+1} = \hat{e}^t - x_j (\hat{b}_j^{t+1} - \hat{b}_j^t)$$

}



<sup>1</sup> Legarra, A., & Misztal, I. (2008). Computing strategies in genome-wide selection. *Journal of dairy science*, 91(1), 360-366.

<sup>2</sup> Xavier, A. (2021). Technical nuances of machine learning. *Crop Breeding and Applied Biotechnology*, 21.

# Coefficients for multivariate model

For updating estimated marker effects we define,  $\hat{\beta}_j^{(t)'} = [\hat{\beta}_{j1}^{(t)} \ \hat{\beta}_{j1}^{(t)} \ \dots \ \hat{\beta}_{jK}^{(t)}]$  to be the vector of estimated marker effects for marker  $j$  and all  $K$  environments,  $\mathbf{Z}_j = \oplus_{k=1}^K \mathbf{z}_{jk}$  to be a matrix containing marker scores at marker  $j$ , and  $\hat{\Sigma}_e^{(t)} = \text{Diag}\{\hat{\sigma}_{e1}^{2(t)}, \hat{\sigma}_{e2}^{2(t)}, \dots, \hat{\sigma}_{ek}^{2(t)}\}$  to be a diagonal matrix of estimated residual variances. Effects for marker  $j$  are initialized with zero and updated as

$$\hat{\beta}_j^{(t+1)} = (\hat{\Sigma}_e^{-1(t)} \mathbf{Z}_j' \mathbf{Z}_j + \hat{\Sigma}_\beta^{-1(t)})^{-1} \mathbf{Z}_j' \hat{\Sigma}_e^{-1(t)} (\mathbf{Z}_j \hat{\beta}_j^{(t)} + \hat{e}^{(t)}), \quad (5)$$

and before moving to the next marker, the residual vector is updated as

$$\hat{e}^{(t+1)} = \hat{e}^{(t)} - \mathbf{Z}_j' (\hat{\beta}_j^{(t+1)} - \hat{\beta}_j^{(t)}). \quad (6)$$

Note that the computation of Kronecker products are not necessary for the multivariate Gauss-Seidel formulation (eq. 5) as long as the residual covariance  $\hat{\Sigma}_e$  is a diagonal matrix.

NO KRONECKER PRODUCTS!!!!

For( j in 1:p ) {

These genetic covariances are the whole key for the MRR model

1<sup>st</sup> solve for beta

$$\begin{bmatrix} \hat{\Sigma}_{\beta}^{11} + \mathbf{z}'_{j(1)}\mathbf{z}_{j(1)}\sigma_{e(1)}^{-2} & \hat{\Sigma}_{\beta}^{12} \\ \hat{\Sigma}_{\beta}^{21} & \hat{\Sigma}_{\beta}^{22} + \mathbf{z}'_{j(2)}\mathbf{z}_{j(2)}\sigma_{e(2)}^{-2} \end{bmatrix} \begin{bmatrix} \hat{\beta}_{j(1)}^{t+1} \\ \hat{\beta}_{j(2)}^{t+1} \end{bmatrix} = \begin{bmatrix} \sigma_{e(1)}^{-2} (\mathbf{z}'_{j(1)}\mathbf{z}_{j(1)}\hat{\beta}_{j(1)}^t + \mathbf{z}'_{j(1)}\hat{e}_1^t) \\ \sigma_{e(2)}^{-2} (\mathbf{z}'_{j(2)}\mathbf{z}_{j(2)}\hat{\beta}_{j(2)}^t + \mathbf{z}'_{j(2)}\hat{e}_2^t) \end{bmatrix}$$

2<sup>nd</sup> update residuals

$$\begin{bmatrix} \hat{e}_{j(1)}^{t+1} \\ \hat{e}_{j(2)}^{t+1} \end{bmatrix} = \begin{bmatrix} \hat{e}_1^t + \mathbf{z}'_{j(1)}(\hat{\beta}_{j(1)}^{t+1} - \hat{\beta}_{j(1)}^t) \\ \hat{e}_2^t + \mathbf{z}'_{j(2)}(\hat{\beta}_{j(2)}^{t+1} - \hat{\beta}_{j(2)}^t) \end{bmatrix}$$

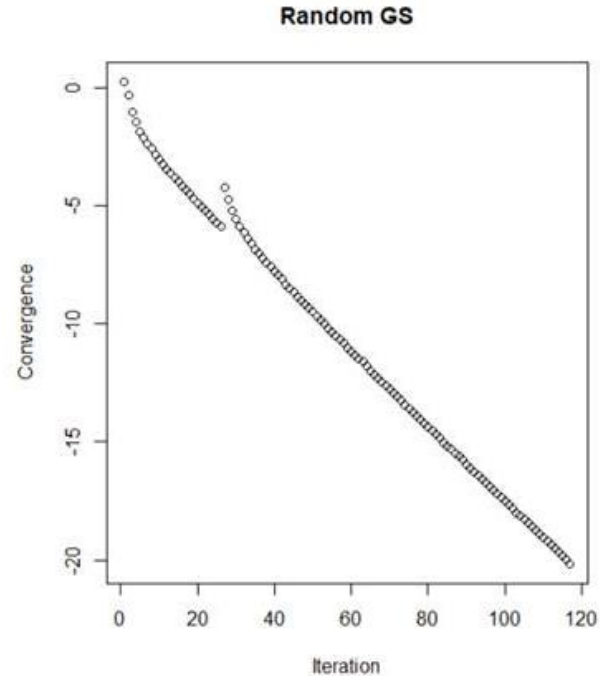
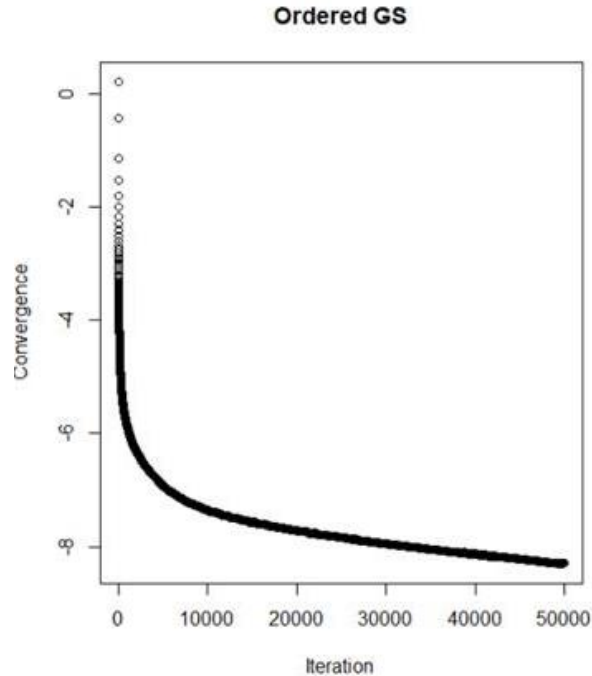
**Color code**

- **Computed only once, before the loop starts (ZpZ)**
- **Computed once every iteration**
- **Computed for each marker in every iteration**

What is in memory?

- |             |                                       |
|-------------|---------------------------------------|
| - Z (n x m) | - ZpZ (m x k)                         |
| - B (m x k) | - $\hat{\Sigma}_{\beta}^{-1}$ (k x k) |
| - E (n x k) | - $\hat{\Sigma}_e^{-1}$ (k)           |

# Side note: Updating markers in random order can speed up convergence



# 1. Introduction

- Rationale and statistical model

# 2. Coefficients

- Univariate
- Multivariate

# 3. Variances

- Univariate
- Multivariate

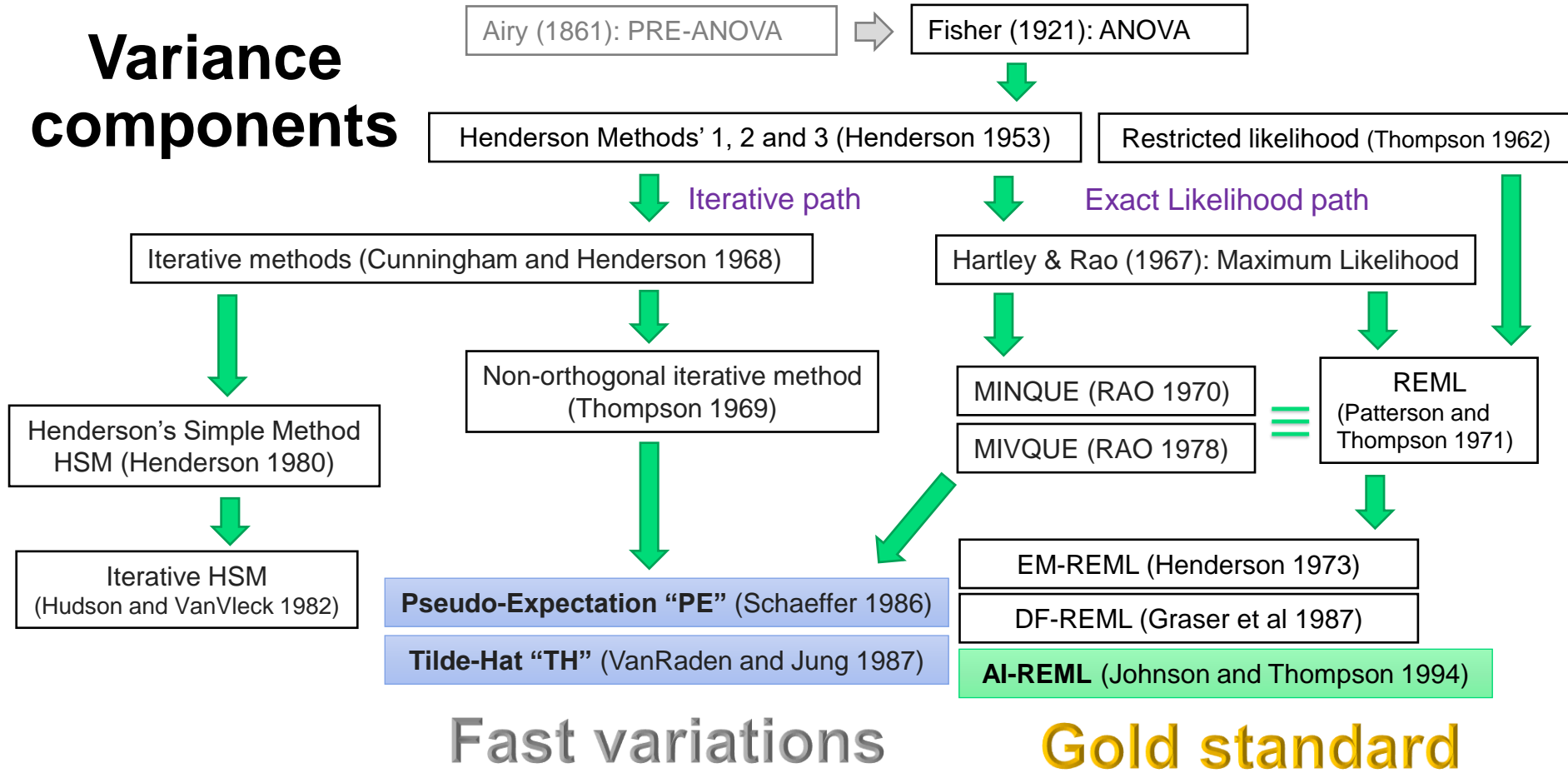
# 4. Simulations

- Study 1: Comparison to REML in small balanced data
- Study 2: Performance in large unbalanced data
- Limitations and other considerations

# 5. Conclusion



# Variance components



# Univariate case: Variance components

- REML

$$\hat{\sigma}_{\beta}^2 = \frac{y'P'V_iPy}{\text{tr}(PV_i)} = \frac{y'S'V^{-1}ZZ'V^{-1}Sy}{\text{tr}(V^{-1}SZZ')} = \frac{\hat{\beta}\hat{\beta}}{\text{tr}(V^{-1}\tilde{Z}'\tilde{Z})}$$

"Let's get rid of this  $V^{-1}$ !"

- Schaffer's (Thompson's) Pseudo-Expectation

$$\hat{\sigma}_{\beta}^2 = \frac{y'S'\cancel{V^{-1}}ZZ'\cancel{V^{-1}}Sy}{\text{tr}(\cancel{V^{-1}}SZZ')} = \frac{\tilde{y}'Z\hat{\beta}}{\text{tr}(\tilde{Z}'\tilde{Z})}$$

"Let's replace this  $V^{-1}$  by something similar, but easier to compute!"

- VanRaden's Tilde-Hat

$$\hat{\sigma}_{\beta}^2 = \frac{y'S'D^{-1}ZZ'V^{-1}Sy}{\text{tr}(D^{-1}SZZ')} = \frac{\tilde{y}\overbrace{D^{-1}Z}^{\tilde{\beta}}\hat{\beta}}{\text{tr}(D^{-1}\tilde{Z}'\tilde{Z})} = \frac{\tilde{\beta}\hat{\beta}}{\text{tr}(D^{-1}\tilde{Z}'\tilde{Z})}$$

All methods yield the same residual variance:

$$\hat{\sigma}_e^2 = \frac{y'e}{n-1}$$

V is a pain to compute

$$V = ZZ'\sigma_{\beta}^2 + I\sigma_e^2$$

$$S = I - (X'X)^{-1}X'; \quad P = V^{-1}S$$

$$P = V^{-1} - V^{-1}(X'V^{-1}X)^{-1}X'V^{-1}$$

$$PX = SX = 0$$

$$Sy = \text{Centralized } y = \tilde{y}$$

$$SZ = \text{Centralized } Z = \tilde{Z}$$

$$D = \text{Diag}(Z'Z\hat{\sigma}_e^{-2} + I\hat{\sigma}_{\beta}^{-2})$$

# Multivariate case: (co)variance components

$$\hat{\sigma}_{\beta(k)}^2 = \frac{\tilde{\beta}_k \hat{\beta}_k}{\text{tr}(\mathbf{D}_k^{-1} \tilde{\mathbf{Z}}_k' \tilde{\mathbf{Z}}_k)} \quad \hat{\sigma}_{\beta(k,k')} = \frac{\tilde{\beta}_k \hat{\beta}_{k'} + \tilde{\beta}_{k'} \hat{\beta}_k}{\text{tr}(\mathbf{D}_k^{-1} \tilde{\mathbf{Z}}_k' \tilde{\mathbf{Z}}_k) + \text{tr}(\mathbf{D}_{k'}^{-1} \tilde{\mathbf{Z}}_{k'}' \tilde{\mathbf{Z}}_{k'})}$$

$$\hat{\sigma}_{e(k)}^2 = \frac{y_k' \hat{e}_k}{n_k - 1}$$

Note: Schaffer's is obtained by assuming  $\mathbf{D} = \mathbf{I}$

**No  $\mathbf{V}$ , No  $\mathbf{C}$ , No LHS,  
No determinants,  
No dense inversions**

## Color code

- Computed only once, before the loop starts (ZpZ)
- Computed once every iteration
- Computed once for PE, and every iteration for TH

What is in memory? -  $\mathbf{Y}$  (n x k) -  $\hat{\Sigma}_{\beta}$  (k x k)

- $\mathbf{Z}$ (n x m)	- $\mathbf{Y}_{\text{tilde}}$ (n x k)	- $\hat{\Sigma}_e$ (k)
- $\mathbf{B}_{\text{hat}}$ (m x k)	- $\mathbf{ZpZ}$ (m x k)	- $\mathbf{N}$ (k)
- $\mathbf{B}_{\text{tilde}}$ (m x k)	- $\mathbf{ZpZ}_{\text{tilde}}$ (m x k)	
- $\mathbf{E}$ (n x k)		

# An intuitive derivation for Schaeffer's method?

The genetic covariance is simply estimated as the cross-prediction between traits A and B normalized by mean squared genotypes (MSX)!!

$$\hat{\sigma}_{\beta(A,B)} = \frac{\overset{\text{Centered phenotype of A}}{(y_A - \mu_A)}' \overset{\text{A predicted from B}}{(Z_A \beta_B)} + \overset{\text{Centered phenotype of B}}{(y_B - \mu_B)}' \overset{\text{B predicted from A}}{(Z_B \beta_A)}}{\text{MSX}_A + \text{MSX}_B}$$

$$*\text{MSX} = \text{Tr}(\tilde{\mathbf{Z}}'\tilde{\mathbf{Z}}) = n \sum_{j=1}^P \hat{\sigma}_{\tilde{\mathbf{Z}}_j}^2$$

# The key parameters from multivariate models

- Genetic variance

$$\hat{\sigma}_{a(k)}^2 = \hat{\sigma}_{\beta(k)}^2 \text{tr}(\mathbf{D}_k^{-1} \tilde{\mathbf{Z}}_k' \tilde{\mathbf{Z}}_k)$$

- Heritability

$$\hat{h}_{(k)}^2 = \frac{\hat{\sigma}_{a(k)}^2}{\hat{\sigma}_{a(k)}^2 + \hat{\sigma}_{e(k)}^2}$$

- Genetic correlations

$$\hat{\rho}_{(k,k')} = \frac{\hat{\sigma}_{\beta(k,k')}}{\sqrt{\hat{\sigma}_{a(k)}^2 \hat{\sigma}_{a(k')}^2}}$$

# 1. Introduction

- Rationale and statistical model

# 2. Coefficients

- Univariate
- Multivariate

# 3. Variances

- Univariate
- Multivariate

# 4. Simulations

- Study 1: Comparison to REML in small balanced data
- Study 2: Performance in large unbalanced data
- Limitations and other considerations

# 5. Conclusion

# Metrics

## 1. Breeding values:

$$\text{Accuracy} = \text{cor}(\text{GEBV}, \text{TBV})$$

## 2. Heritability ( $h^2$ ) and genetic correlations ( $\rho$ ):

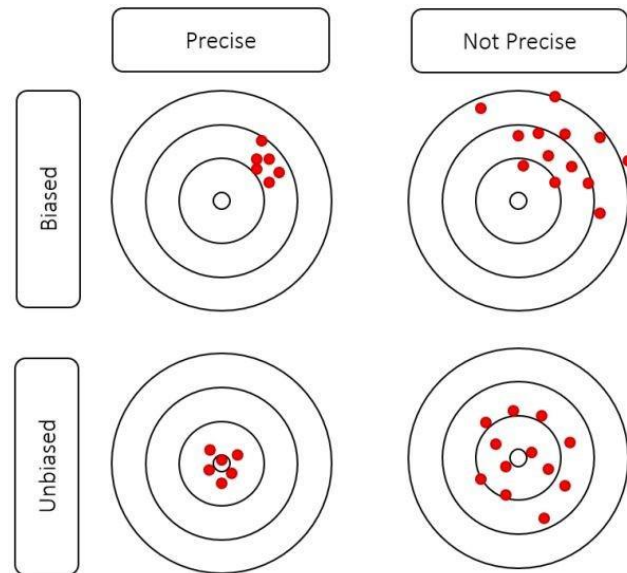
$$\text{Bias} = E(\hat{\theta} - \theta)$$

$$\text{Precision} = \text{SD}(\hat{\theta} - \theta)$$



## 3. Computation efficiency:

Elapsed time to fit the model



[Picture source](#)

# Study 1

- Wheat data (CYMMIT)
- 599 Individuals
- 1299 Markers
- Scenario: 10 environments, all individuals observed in all locations
- Methods: REML, PEGS, THGS, Univariate

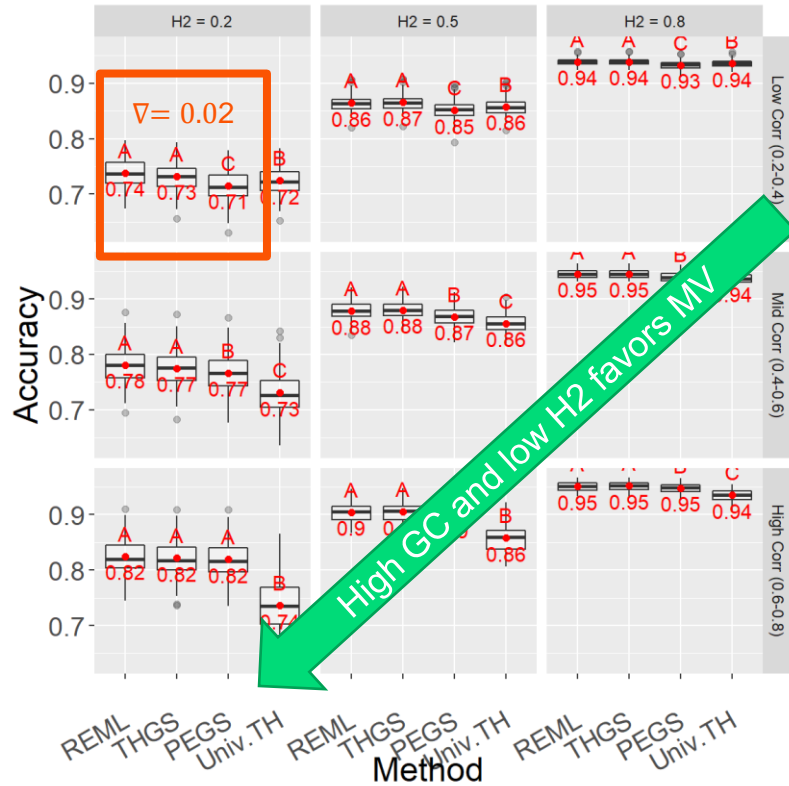


# Elapsed time

Method	Time in minutes (SD)	
REML	<b>256.9 (60.57)</b>	= 4 hours and 17 minutes
PEGS	0.27 (0.02)	= 16 seconds
THGS	0.27 (0.02)	
Univariate	0.23 (0.03)	= 13 seconds

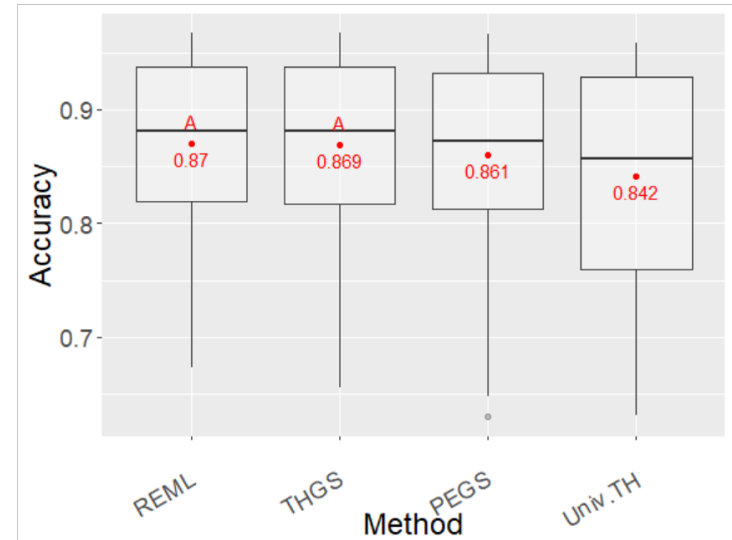
Wheat dataset: 10 traits, 599 individuals, 1299 markers  
(available in the BGLR package)

# Accuracy of breeding values

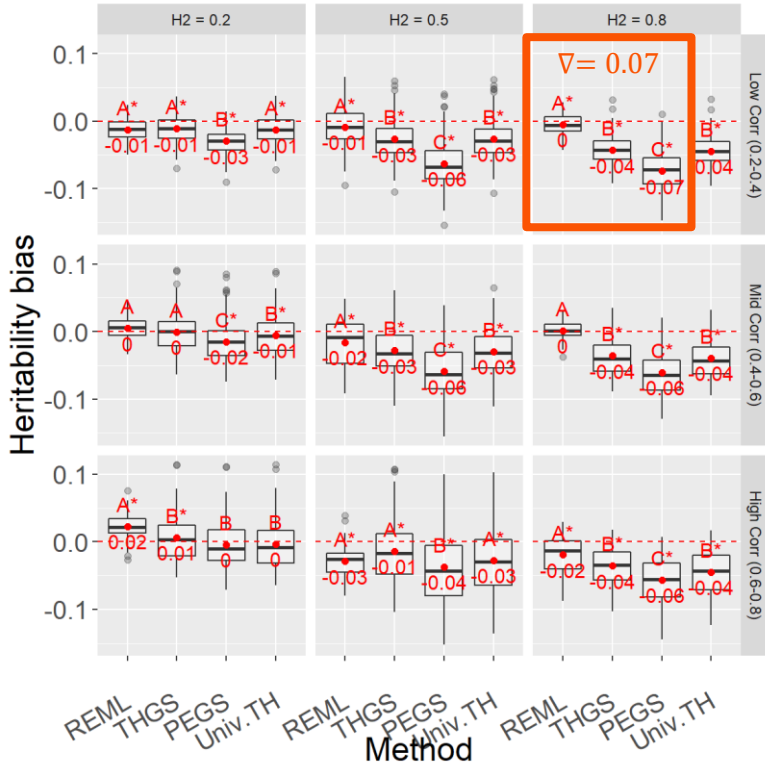


$$Acc = \text{cor}(\text{GEBV}, \text{TBV})$$

(Higher is better)

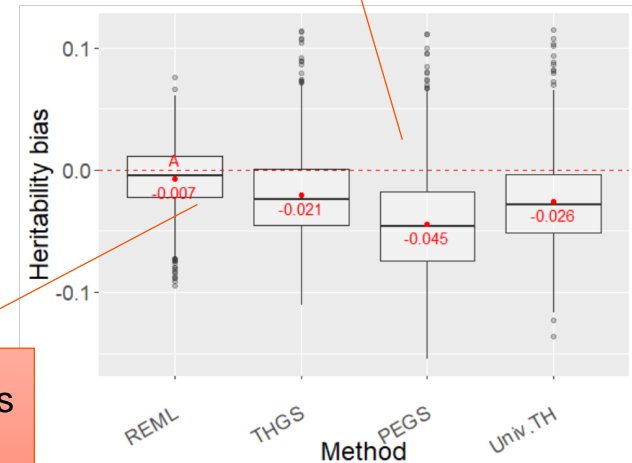


# Bias of heritability estimates



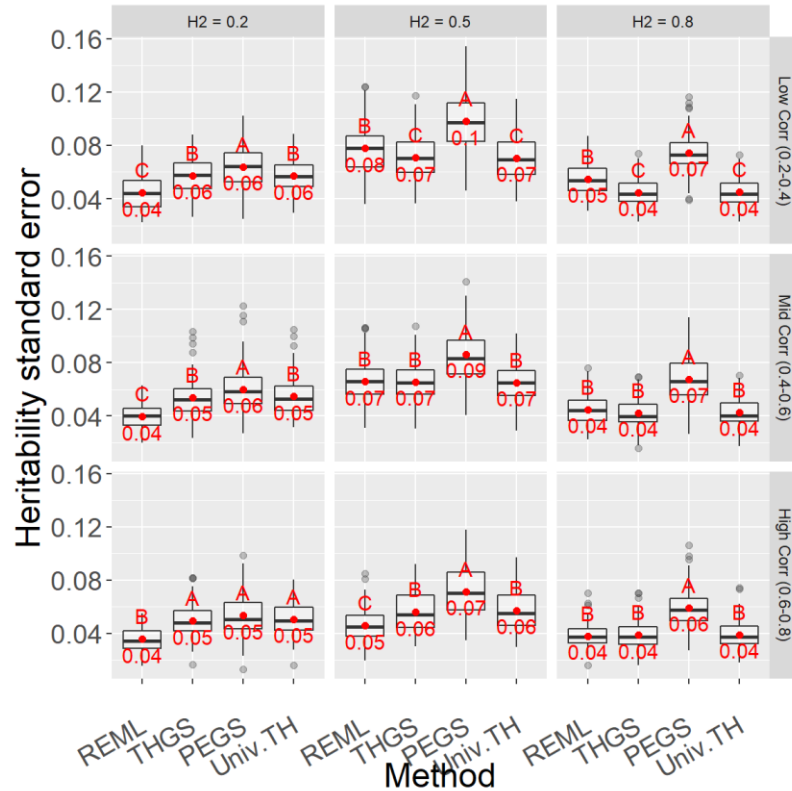
**Bias  $h^2 = E(\hat{h}^2 - h^2)$**   
(Closer to zero is better)

PEGS underestimated  $h^2$   
when true  $h^2$  was mid-high



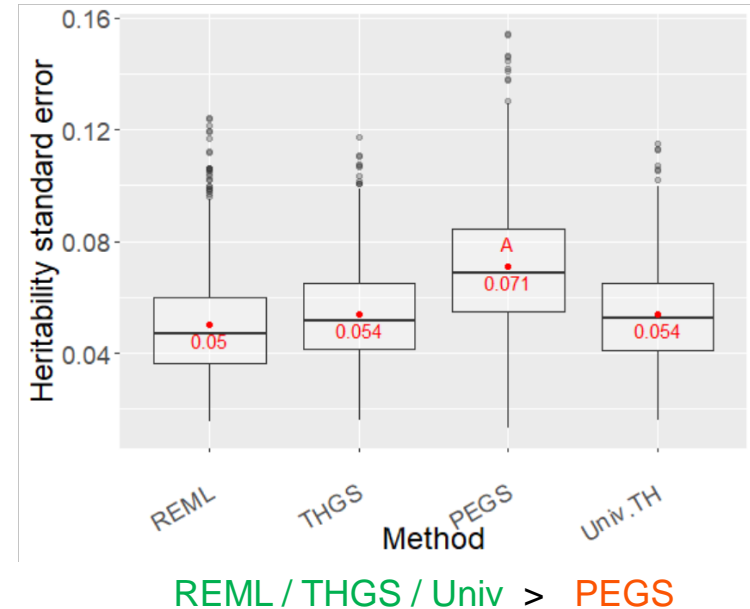
REML also shows large variation...

# Precision of heritability estimates



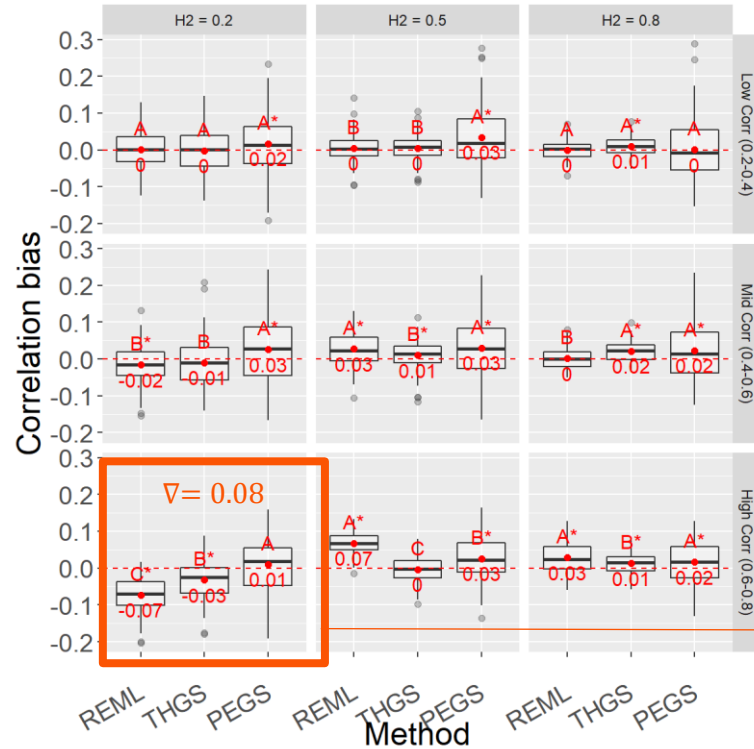
$$\text{Prec } h^2 = \text{SE}(\hat{h}^2 - h^2)$$

(Lower is better)



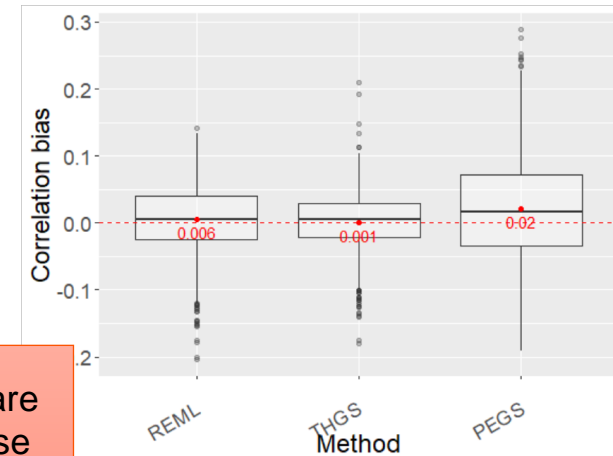
REML / THGS / Univ > PEGS

# Bias of genetic correlation estimates



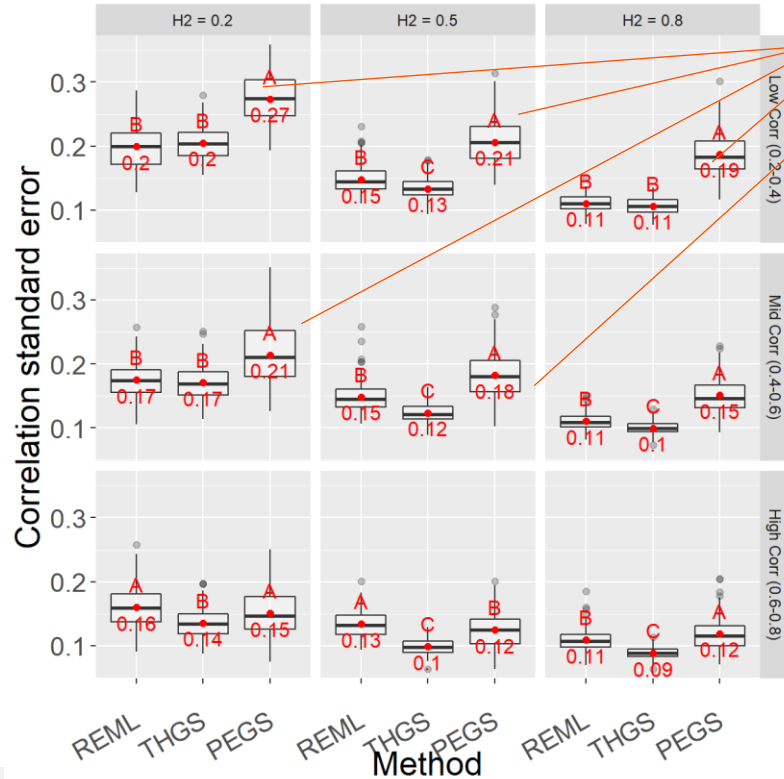
$$\text{Bias } \rho = E(\hat{\rho} - \rho)$$

(Closer to zero is better)



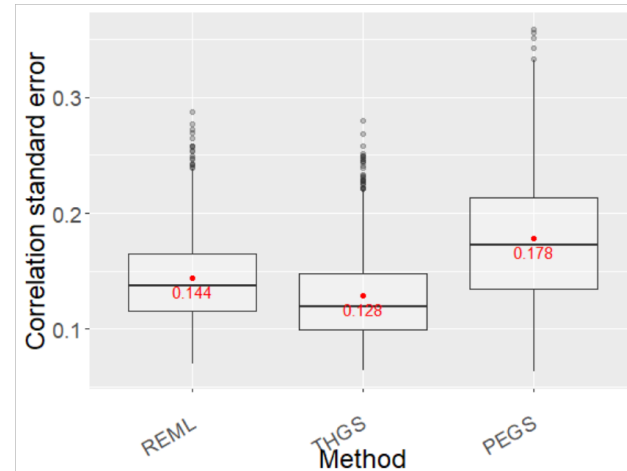
Differences are large because REML is doing a poor job

# Precision of genetic correlation estimates



PEGS has a hard time to estimate correlations when heritability is low, possibly because it underestimates Genetic Variances

Precision  $\rho = SE(\hat{\rho} - \rho)$   
(Lower is better)



THGS > REML > PEGS

# Summary of the smaller & balanced (wheat) dataset

Method	Accuracy	Bias H2	Precision H2	Bias GC	Precision GC
REML	<b>0.88 (0.01)</b>	-0.00 (0.03)	0.04 (0.02)	0.01 (0.05)	0.15 (0.03)
PEGS	0.87 (0.02)	-0.03 (0.02)	0.04 (0.01)	0.01 (0.08)	<b>0.18*</b> (0.04)
THGS	<b>0.88 (0.01)</b>	-0.01 (0.01)	<b>0.03 (0.01)</b>	<b>-0.01 (0.04)</b>	<b>0.13 (0.02)</b>
Univariate	0.85 (0.03)	-0.01 (0.01)	<b>0.03 (0.01)</b>	-	-

\* PEGS correlations were less precise than THGS, but not statistically different than REML in small balanced datasets

# Study 2

- Soybean data (SoyNAM)
- 5000 Individuals
- 4300 Markers
- Scenario: 10 environments, no overlapping individuals
  - **Each individual is observed in a single environment!**
- Methods: PEGS, THGS, Univariate



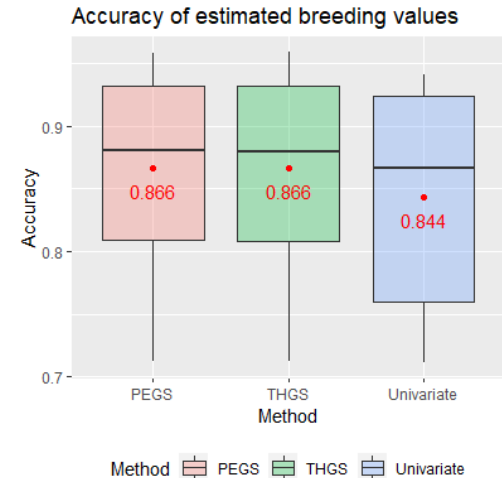
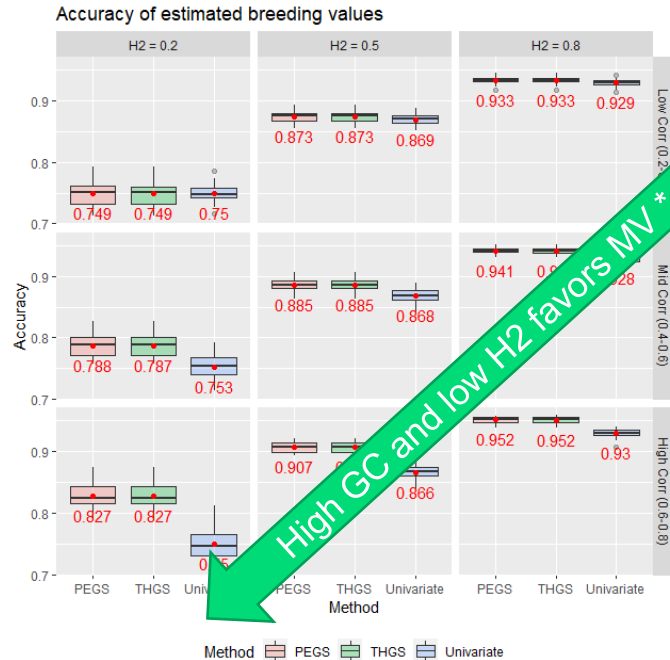
# Elapsed time

No. of environments		PEGS	THGS	Univariate-THGS
10		0.7 (0.2)	0.7 (0.2)	0.2 (0.0)
50		12 (3)	12 (3)	1.0 (0.1)
100		43 (13)	44 (14)	2.0 (0.3)
200	~3h	168 (48)	165 (44)	4.0 (0.4)
400	~10h	568 (47)	560 (53)	8.0 (1.9)
500	~14h	807 (39)	832 (49)	10.0 (0.6)

(Time in minutes)

# Accuracy of breeding values

$$Acc = \text{cor}(\text{GEBV}, \text{TBV})$$

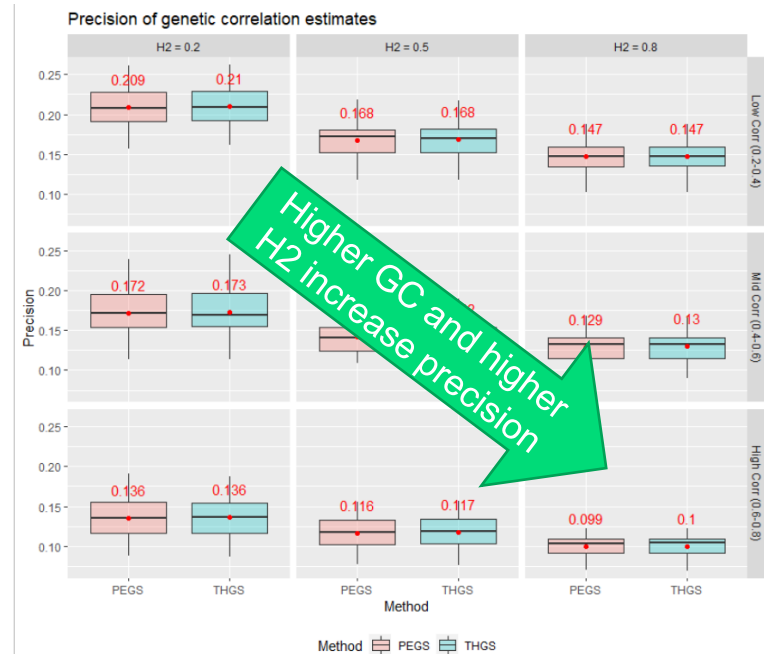
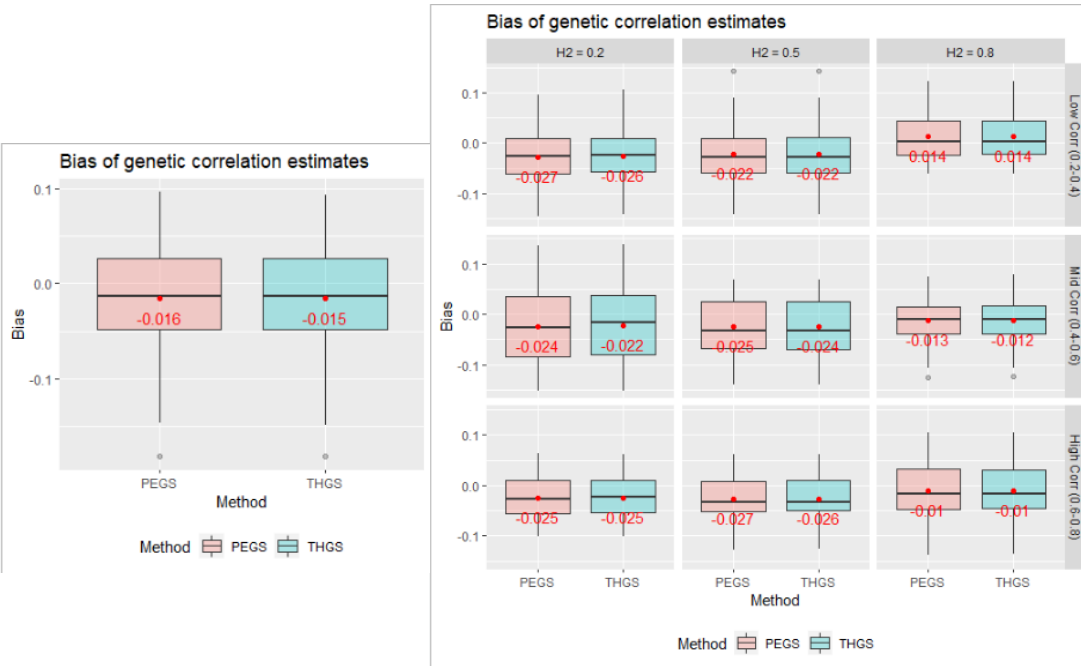


\* Same trend observed in wheat

# Bias of genetic correlation estimates

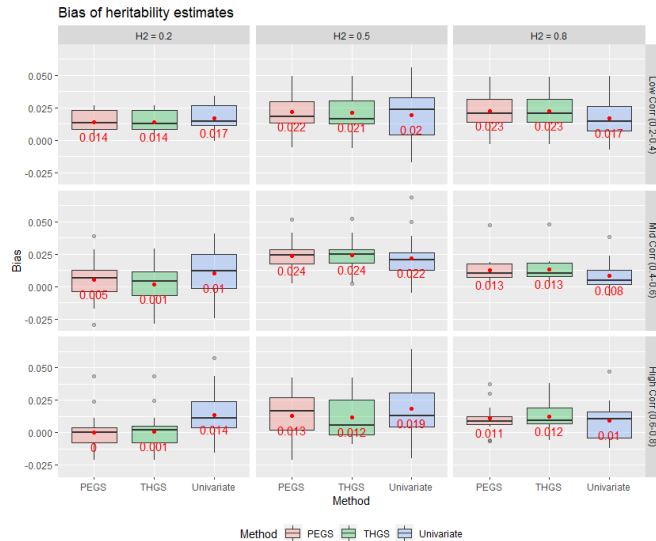
$$\text{Bias } \rho = E(\hat{\rho} - \rho)$$

$$\text{Precision } \rho = SE(\hat{\rho} - \rho)$$

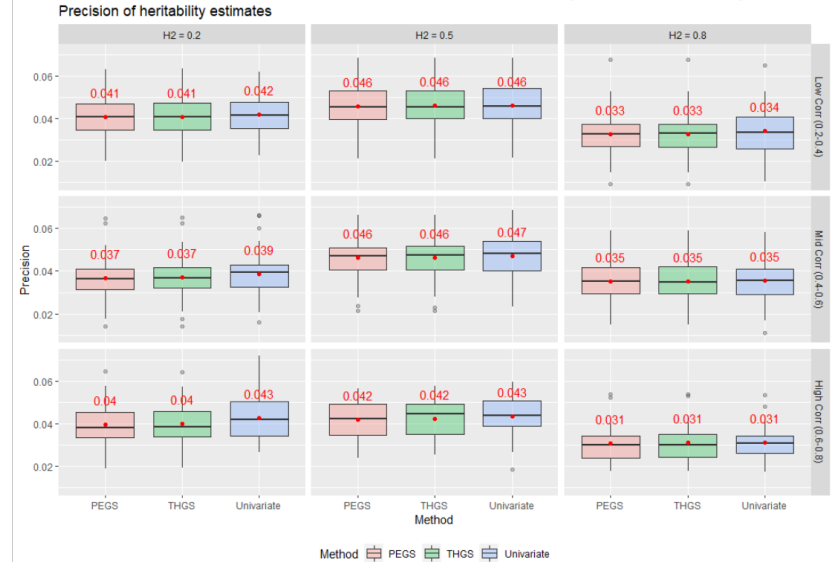


# Bias of heritability estimates

$$\text{Bias } h^2 = E(\hat{h}^2 - h^2)$$



$$\text{Precision } h^2 = SE(\hat{h}^2 - h^2)$$



All roughly the same ~ bias 0.01, S.E. 0.04

# Summary in smaller balanced dataset (wheat)

Method	Time (in min.)	Accuracy	Bias H2	Precision H2	Bias GC	Precision GC
REML	256.90 (60.57)	0.88 (0.01)	-0.00 (0.03)	0.04 (0.02)	0.01 (0.05)	0.15 (0.03)
PEGS	0.27 (0.02)	0.87 (0.02)	-0.03 (0.02)	0.04 (0.01)	0.01 (0.08)	0.18 (0.04)
THGS	0.27 (0.02)	0.88 (0.01)	-0.01 (0.01)	0.03 (0.01)	-0.01 (0.04)	0.13 (0.02)
Univariate	0.23 (0.03)	0.85 (0.03)	-0.01 (0.01)	0.03 (0.01)	-	-

THGS  $\geq$  REML  $\geq$  PEGS  $>$  Univ

# Summary in larger unbalanced dataset (soy)

Method	Accuracy	Bias H2	Prec. H2	Bias GC	Prec. GC
PEGS	0.87 (0.01)	-0.01 (0.01)	0.04 (0.01)	-0.02 (0.06)	0.14 (0.02)
THGS	0.87 (0.01)	-0.01 (0.01)	0.04 (0.01)	-0.02 (0.06)	0.14 (0.02)
Univariate	0.85 (0.02)	-0.02 (0.02)	0.04 (0.01)	-	-

PEGS  $\cong$  THGS  $>$  Univ

# Limitations and other considerations

- **More fixed effects?** The absorption of fixed effects beyond the intersect can create a large computational burden. But it is OK to work with pre-adjusted phenotypes like BLUEs, BLUPs and deregressed BLUPs<sup>1</sup>.
- **Correlated residuals:** Modeling residual covariances may offset most saving in computation time because of the need for  $n \times n$  Kronecker products.
- **Kernels & SVD:** When  $P \gg N$ , Gauss-Seidel may be costly. When feasible, a solution comes from regress Eigenvectors<sup>2</sup> instead ( $Z=UDV$ , solve the MRR using  $Z^*=UD$ , back solve coefficients  $\beta = \beta^*V$ ).
- **Bending**<sup>3</sup>: The covariance  $\hat{\Sigma}_\beta$  may not be invertible with too many correlated traits. One may need to shrink the covariance until  $\hat{\Sigma}_\beta$  can be inverted. Alternatively, use of simpler covariances: CS and XFA.
- **Balanced data:** REML can be efficiently computed when all phenotypes are collected in all individuals using *canonical transformation*<sup>4</sup> or *kernel diagonalization via eigendecomposition*<sup>5</sup>

1 Garrick et al (2009). Deregressing estimated breeding values and weighting information for genomic regression analyses. Genetics Selection Evolution, 41(1), 1-8.

2 Ødegård et al (2018). Large-scale genomic prediction using singular value decomposition of the genotype matrix. Genetics Selection Evolution, 50(1), 1-12.

3 Jorjani et al (2003). A simple method for weighted bending of genetic (co) variance matrices. Journal of dairy science, 86(2), 677-679.

4 Meyer, K. (1985). Maximum likelihood estimation of variance components for a multivariate mixed model with equal design matrices. Biometrics, 153-165.

5 Lee and Van der Werf (2016). MTG2: an efficient algorithm for multivariate linear mixed model analysis based on genomic information. Bioinformatics, 32(9), 1420-1422.

## 1. Introduction

- Rationale and statistical model

## 2. Coefficients

- Univariate
- Multivariate

## 3. Variances

- Univariate
- Multivariate

## 4. Simulations

- Study 1: Comparison to REML in small balanced data
- Study 2: Performance in large unbalanced data
- Limitations and other considerations

## 5. Conclusion

# Thank you for your attention!

## Remarks:

- 1) Multivariate models are valuable, but these have been computationally unfeasible
- 2) Efficient estimation of coefficients (RGS) and variances (PE/TH) enable large MRR
- 3) THGS & PEGS have some limitations but are suitable replacements to REML

## Questions??

***Alencar Xavier***

[Alencar.Xavier@Corteva.com](mailto:Alencar.Xavier@Corteva.com)