

# Genetic association and heritability estimation in structured populations

Alejandro Ochoa

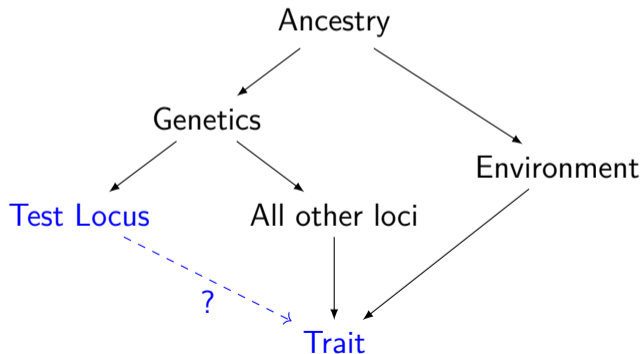
StatGen, Biostatistics & Bioinformatics — Duke University

2023-11-18 — CBB retreat

Part 1: Cryptic relatedness matters a lot in genetic association studies ... or: how I learned to love Linear Mixed-effects Models

# Association studies are hard

- ▶ Millions of tests
- ▶ Polygenicity (many causal variants)
- ▶ Incorrect assumptions: independence / additivity
- ▶ **Confounders**



# PCA vs LMM in association

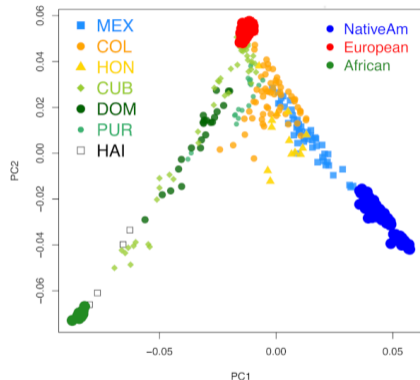
Principal Components Analysis (PCA)  
and Linear Mixed-effects Model (LMM):

$$\text{PCA : } \mathbf{y} = \mathbf{1}\alpha + \mathbf{x}_i\beta + \mathbf{U}_d\gamma_d + \epsilon,$$

$$\text{LMM : } \mathbf{y} = \mathbf{1}\alpha + \mathbf{x}_i\beta + \mathbf{s} + \epsilon.$$

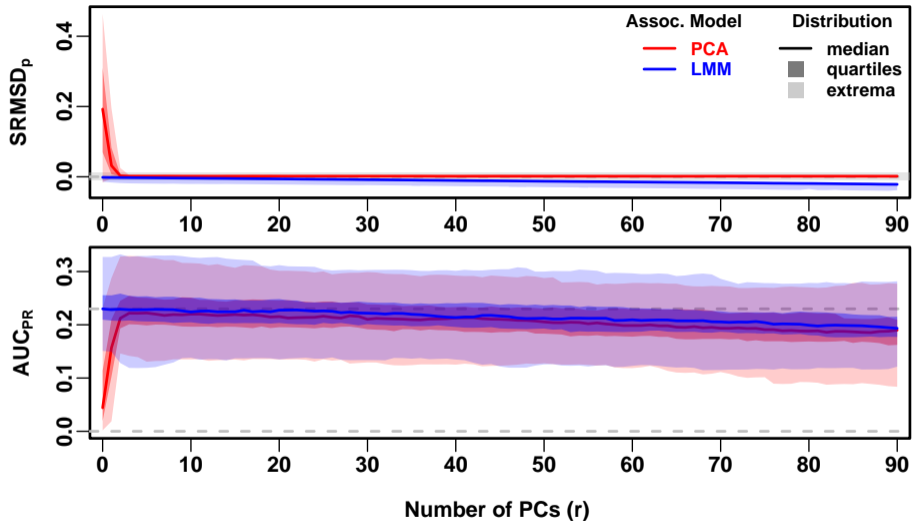
$\mathbf{U}_d$  are top  $d$  eigenvectors of kinship matrix  $\Phi$ .  
 $\mathbf{s} \sim \text{Normal}(\mathbf{0}, \sigma_G^2\Phi)$ .

- ▶ PCA is faster but low-dimensional
- ▶ LMM is slower but can model families

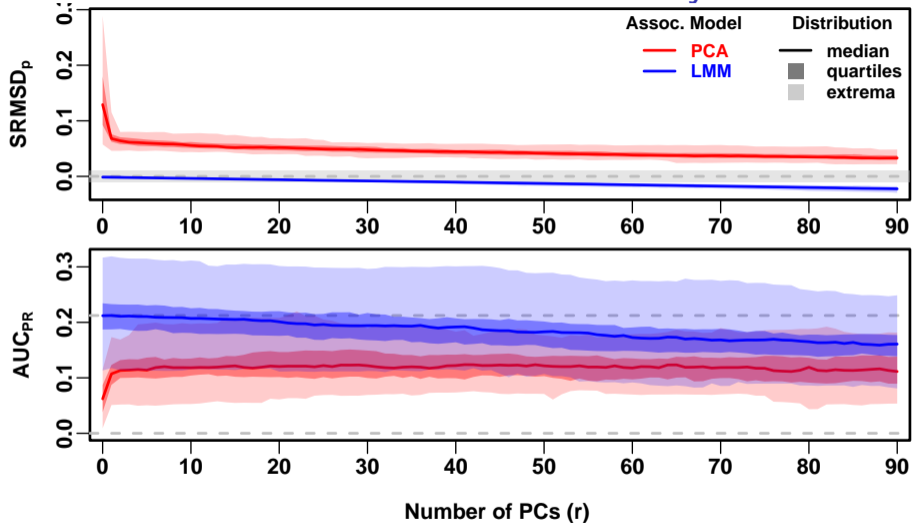


Moreno-Estrada *et al.* (2013)

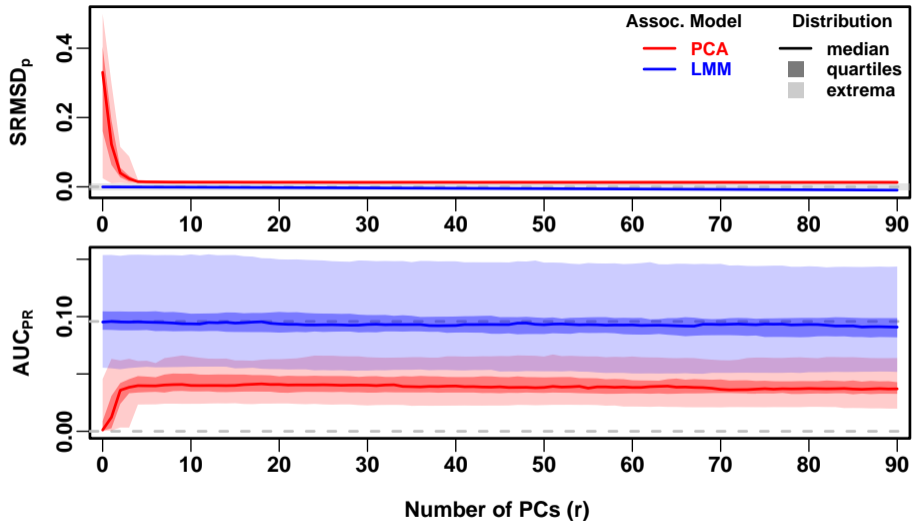
# PCA < LMM in association: simulated admixture



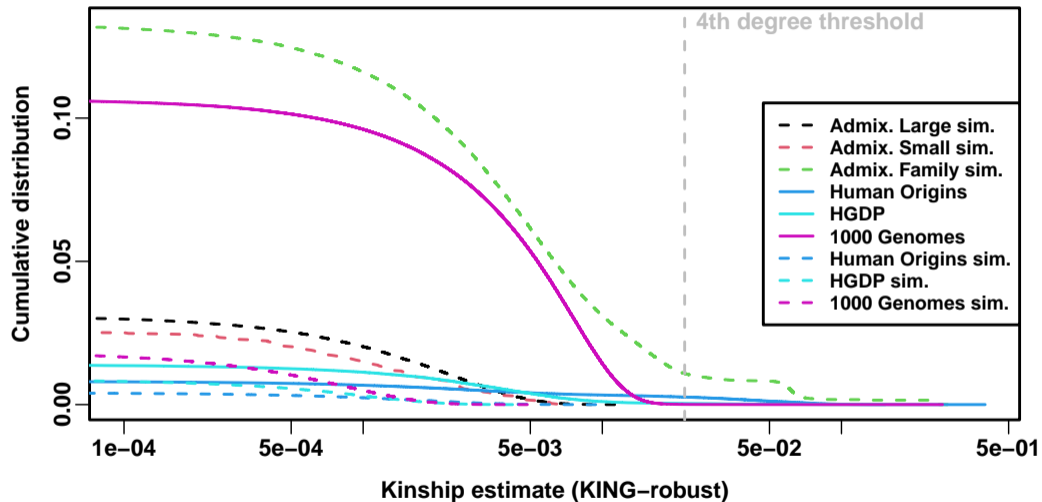
# PCA < LMM in association: simulated family structure



# PCA < LMM in association: 1000 Genomes

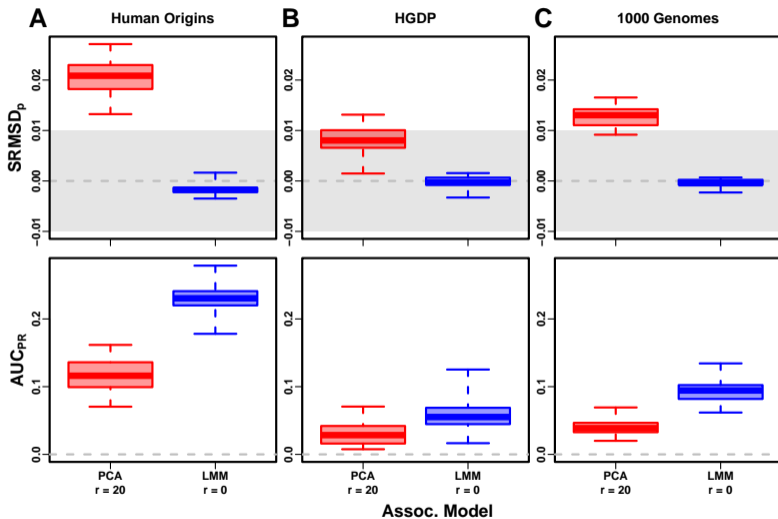


# Numerous distant relatives in real datasets



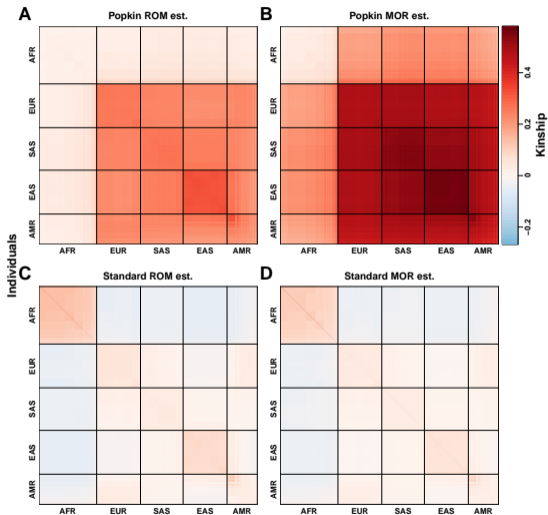


# Numerous distant relatives in real datasets explain $\text{PCA} < \text{LMM}$



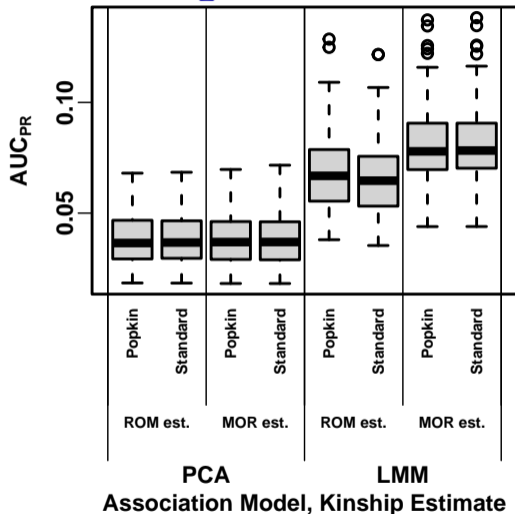
Part 2: Kinship bias carries over to heritability, but not association!

# Standard kinship estimator is severely biased



- ▶ ROM: Ratio of Means
  - ▶ Upweights common variants
  - ▶ Behaves well statistically
- ▶ MOR: Mean of Ratios
  - ▶ Upweights rare variants
  - ▶ Introduces additional bias

# Kinship bias does not affect genetic associations



# Kinship bias does not affect genetic associations

Linear algebra proof!

Transforming true to biased kinship matrices:

$\Phi$  : Unbiased kinship matrix,

$\Phi'$  : Biased kinship matrix,

$$\Phi' = \frac{1}{1 - \bar{\varphi}} \mathbf{C} \Phi \mathbf{C},$$

$$\mathbf{C} = \mathbf{I} - \frac{1}{n} \mathbf{1}\mathbf{1}^\top : \text{Centering matrix.}$$

# Kinship bias does not affect genetic associations

Association test is a regression with correlated residuals:

Linear algebra proof!

Transforming true to biased kinship matrices:

$\Phi$  : Unbiased kinship matrix,

$\Phi'$  : Biased kinship matrix,

$$\Phi' = \frac{1}{1 - \bar{\varphi}} \mathbf{C} \Phi \mathbf{C},$$

$$\mathbf{C} = \mathbf{I} - \frac{1}{n} \mathbf{1}\mathbf{1}^\top : \text{Centering matrix.}$$

$$\mathbf{y} = \mathbf{1}\alpha + \mathbf{x}_i\beta_i + \mathbf{s} + \epsilon,$$

$$\mathbf{s} \sim \text{Normal}(\mathbf{0}, 2\sigma_G^2 \Phi),$$

$$\epsilon \sim \text{Normal}(\mathbf{0}, \sigma_E^2 \mathbf{I}).$$

# Kinship bias does not affect genetic associations

Linear algebra proof!

Transforming true to biased kinship matrices:

$\Phi$  : Unbiased kinship matrix,

$\Phi'$  : Biased kinship matrix,

$$\Phi' = \frac{1}{1 - \bar{\varphi}} \mathbf{C} \Phi \mathbf{C},$$

$$\mathbf{C} = \mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^\top : \text{Centering matrix.}$$

Association test is a regression with correlated residuals:

$$\mathbf{y} = \mathbf{1} \alpha + \mathbf{x}_i \beta_i + \mathbf{s} + \epsilon,$$

$$\mathbf{s} \sim \text{Normal}(\mathbf{0}, 2\sigma_G^2 \Phi),$$

$$\epsilon \sim \text{Normal}(\mathbf{0}, \sigma_E^2 \mathbf{I}).$$

Kinship bias compensated by intercept!

$$\mathbf{s}' = \mathbf{C} \mathbf{s} \sim \text{Normal}(\mathbf{0}, 2\sigma_G'^2 \Phi'),$$

$$\sigma_G'^2 = (1 - \bar{\varphi}) \sigma_G^2,$$

$$\mathbf{s}' = \mathbf{s} - \mathbf{1} \bar{s},$$

$$\alpha' = \alpha + \bar{s}$$

# Kinship bias affects heritability estimation

Model:

$$\mathbf{y} = \mathbf{1}\alpha + \mathbf{s} + \epsilon,$$

$$\mathbf{s} + \epsilon \sim \text{Normal}(\mathbf{0}, 2\sigma_G^2\Phi + \sigma_E^2\mathbf{I}).$$

Heritability definition:

$$h^2 = \frac{\sigma_G^2}{\sigma_G^2 + \sigma_E^2}.$$

Variance is estimated with bias:

$$\sigma_G^{2'} = (1 - \bar{\varphi})\sigma_G^2.$$



# Kinship bias affects heritability estimation

Model:

$$\mathbf{y} = \mathbf{1}\alpha + \mathbf{s} + \epsilon,$$
$$\mathbf{s} + \epsilon \sim \text{Normal}(\mathbf{0}, 2\sigma_G^2\Phi + \sigma_E^2\mathbf{I}).$$

Heritability definition:

$$h^2 = \frac{\sigma_G^2}{\sigma_G^2 + \sigma_E^2}.$$

Variance is estimated with bias:

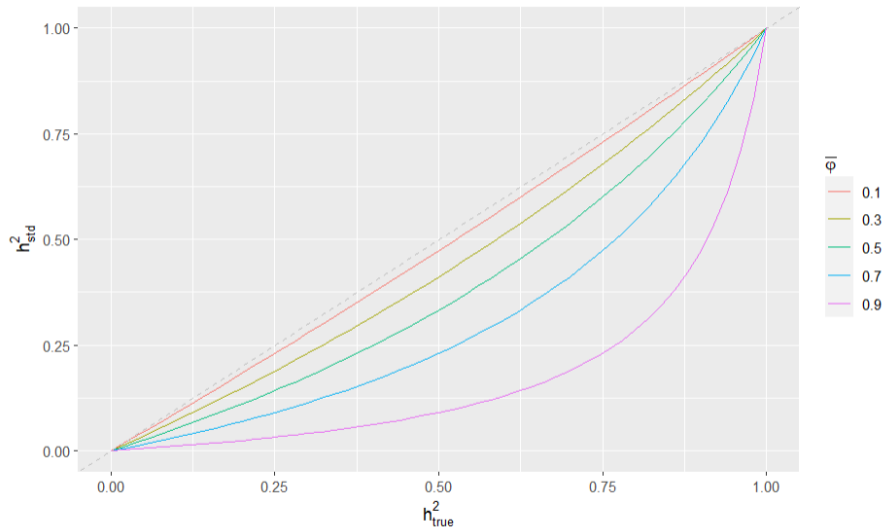
$$\sigma_G^{2'} = (1 - \bar{\varphi})\sigma_G^2.$$

Heritability is estimated with bias that depends on mean kinship  $\bar{\varphi}$  and the true heritability  $h^2$ :

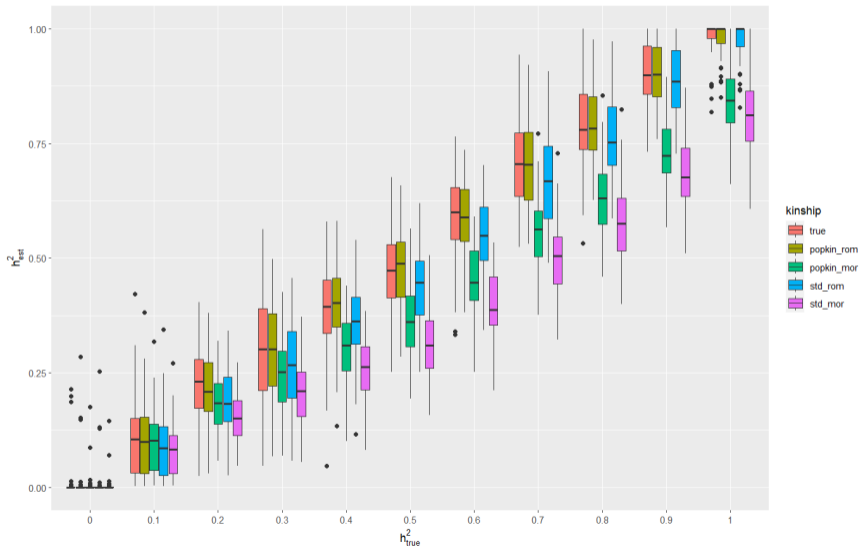
$$h^{2'} = \frac{\sigma_G^{2'}}{\sigma_G^{2'} + \sigma_E^{2'}}$$
$$= h^2 \frac{1 - \bar{\varphi}}{1 - \bar{\varphi}h^2}.$$

# Kinship bias affects heritability estimation

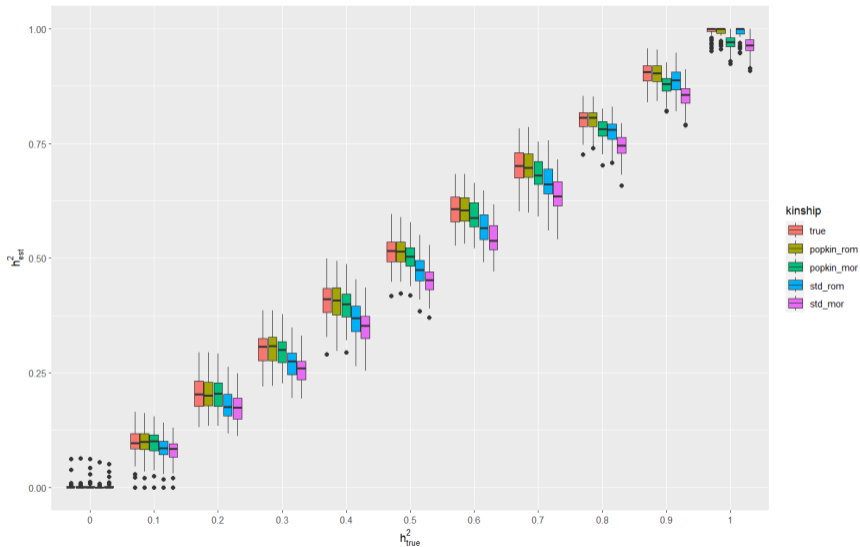
$$h^{2'} = h^2 \frac{1 - \bar{\phi}}{1 - \bar{\phi}h^2}$$



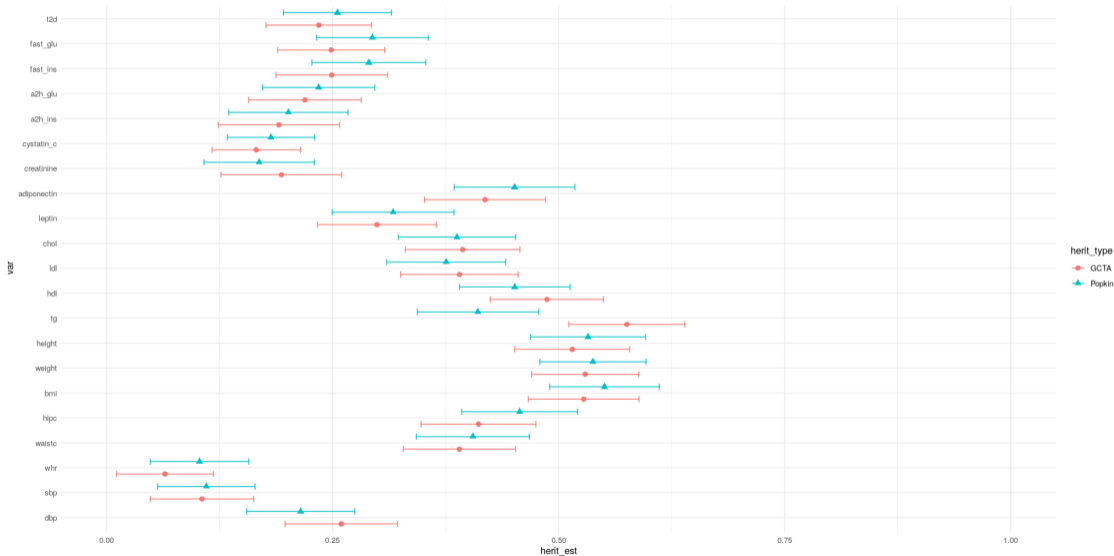
# Simulated admixture (high mean kinship)



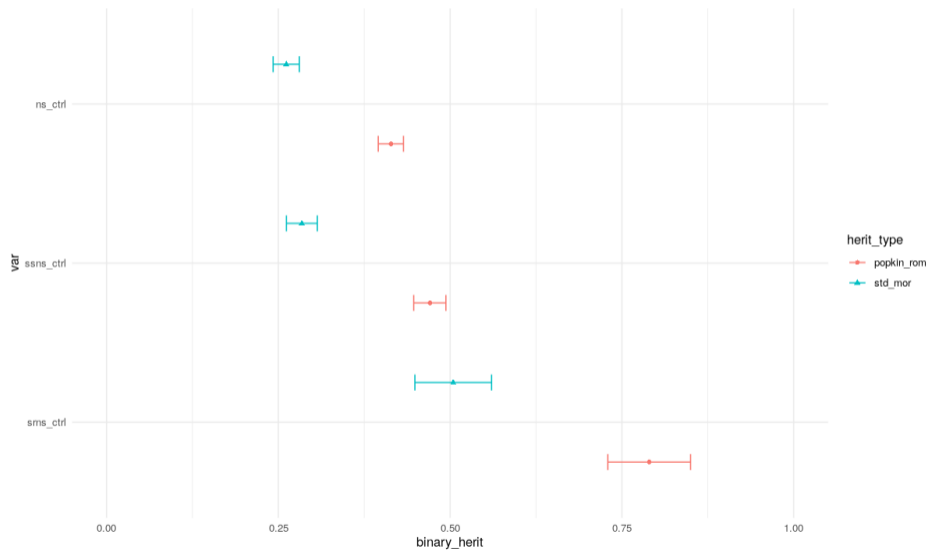
# Simulated family structure (reduced mean kinship)



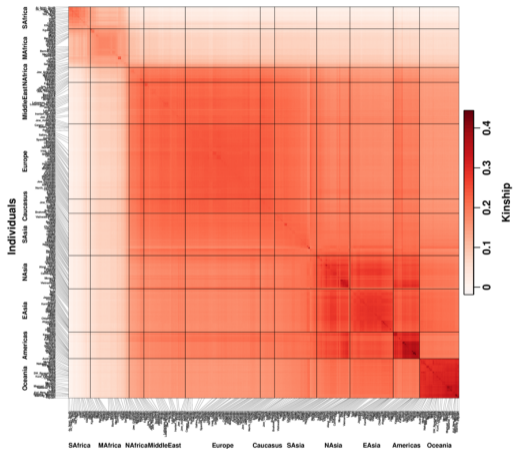
# San Antonio Family Study: Type 2 Diabetes (low mean kinship)



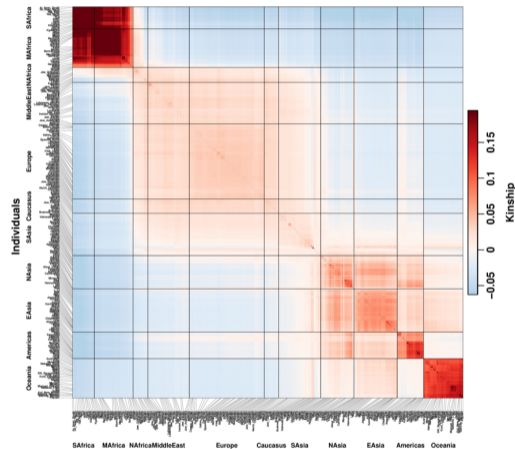
# Nephrotic Syndrome: multiethnic cohort (high mean kinship)



# Unbiased kinship estimates: new models, opportunities



New "popkin"  
kinship estimator



Biased "standard"  
kinship estimator

# Acknowledgments

## Ochoa Lab

Tiffany Tu

RP Pornmongkolsuk

**Zhuoran Hou**

**Yiqi Yao**

Amika Sood

Danielle Mensah

## Princeton University

John D. Storey

## Duke University

Rasheed Gbadegesin

Kouros Owzar

Beth Hauser

Yi-Ju Li

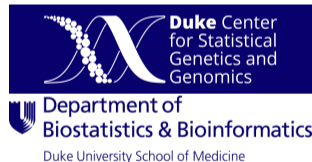
Andrew Allen

Amy Goldberg

## Funding

NIH

Whitehead Scholars



DrAlexOchoa@genomic.social

🏠 [ochoalab.github.io](https://ochoalab.github.io)

✉️ [alejandro.ochoa@duke.edu](mailto:alejandro.ochoa@duke.edu)



# New kinship/GRM estimator

Kinship model for neutral genotypes  $x_{ij} \in \{0, 1, 2\}$ :

$$E[x_{ij}] = 2p_i, \quad \text{Cov}(x_{ij}, x_{ik}) = 4p_i(1 - p_i)\varphi_{jk}.$$

Standard estimator is **biased**:

$$\hat{p}_i = \frac{1}{2n} \sum_{j=1}^n x_{ij}, \quad \hat{\varphi}_{jk}^{\text{std}} = \frac{1}{m} \sum_{i=1}^m \frac{(x_{ij} - 2\hat{p}_i)(x_{ik} - 2\hat{p}_i)}{4\hat{p}_i(1 - \hat{p}_i)} \approx \frac{\varphi_{jk} - \bar{\varphi}_j - \bar{\varphi}_k + \bar{\varphi}}{1 - \bar{\varphi}}.$$

**popkin**: first unbiased kinship estimator! R package

$$A_{jk} = \frac{1}{m} \sum_{i=1}^m (x_{ij} - 1)(x_{ik} - 1) - 1, \quad \hat{\varphi}_{jk}^{\text{new}} = 1 - \frac{A_{jk}}{\hat{A}_{\min}} \xrightarrow[m \rightarrow \infty]{\text{a.s.}} \varphi_{jk}.$$

Ochoa and Storey (2021) doi:10.1371/journal.pgen.1009241

