

Relatedness and differentiation in arbitrary population structures

Alejandro Ochoa, John D. Storey Lab, Princeton University

🐦 DrAlexOchoa

🏠 viiia.org/research/

✉️ ochoa@princeton.edu

My research areas / Contributions

1. Stratified False Discovery Rate (FDR)

- ▶ Finding: per-stratum local FDR maximizes power controlling overall FDR
- ▶ Improved power in protein domain prediction
- ▶ Identified protein classes with problematic statistics

My research areas / Contributions

1. Stratified False Discovery Rate (FDR)

- ▶ Finding: per-stratum local FDR maximizes power controlling overall FDR
- ▶ Improved power in protein domain prediction
- ▶ Identified protein classes with problematic statistics

2. Genetic and other factors controlling the placebo response

- ▶ Collaboration with Otsuka Pharmaceutical
- ▶ Mixed-effects modeling of longitudinal response in drug trials for Schizophrenia, Bipolar Disorder, and Major Depressive Disorder
- ▶ Genome-wide association study for individual-specific placebo response

My research areas / Contributions

1. Stratified False Discovery Rate (FDR)

- ▶ Finding: per-stratum local FDR maximizes power controlling overall FDR
- ▶ Improved power in protein domain prediction
- ▶ Identified protein classes with problematic statistics

2. Genetic and other factors controlling the placebo response

- ▶ Collaboration with Otsuka Pharmaceutical
- ▶ Mixed-effects modeling of longitudinal response in drug trials for Schizophrenia, Bipolar Disorder, and Major Depressive Disorder
- ▶ Genome-wide association study for individual-specific placebo response

3. Kinship and F_{ST} for arbitrary population structures

- ▶ Motivation: world-wide human population structure
- ▶ Generalized definitions and models
- ▶ Novel bias calculations validated by simulations
- ▶ Novel non-parametric estimator with greatly improved accuracy

Why study relatedness?

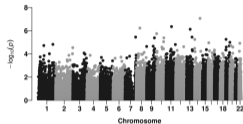


Human genetics
is fascinating!

Why study relatedness?



Human genetics
is fascinating!



Pop. structure
confounds
association
studies (GWAS)

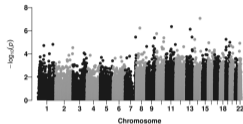
Why study relatedness?



Human genetics
is fascinating!



Heritability of
complex traits

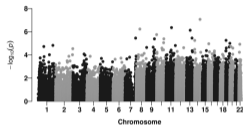


Pop. structure
confounds
association
studies (GWAS)

Why study relatedness?



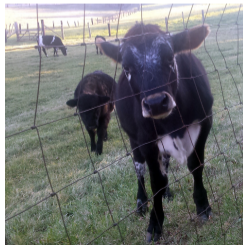
Human genetics
is fascinating!



Pop. structure
confounds
association
studies (GWAS)

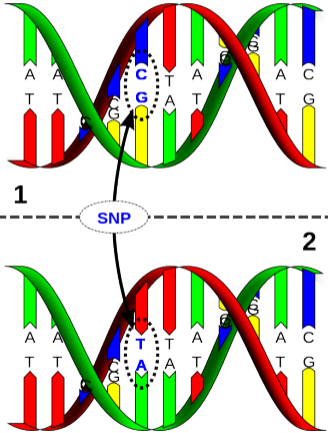


Heritability of
complex traits

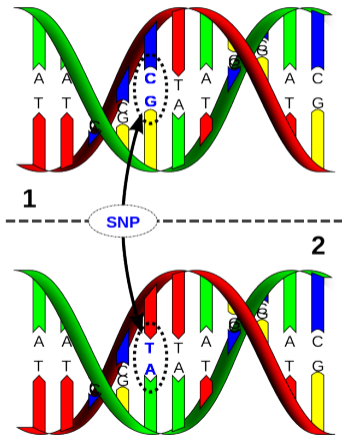


Animal and plant
breeding

Single Nucleotide Polymorphism (SNP) data



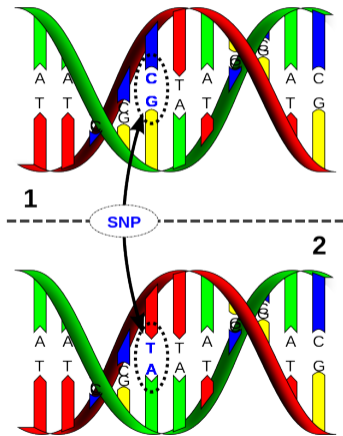
Single Nucleotide Polymorphism (SNP) data



⇒

Genotype	x_{ij}
CC	0
CT	1
TT	2

Single Nucleotide Polymorphism (SNP) data



⇒

Genotype	x_{ij}
CC	0
CT	1
TT	2

⇒

	Individuals						
Loci	0	2	2	1	1	0	1
	0	2	1	0	1		
	2	...					

X

Hardy-Weinberg Equilibrium (HWE): Binomial draws

x_{ij} = genotype at locus i for individual j .

p_i^T = frequency of reference allele at locus i , (ancestral) population T .

Hardy-Weinberg Equilibrium (HWE): Binomial draws

x_{ij} = genotype at locus i for individual j .

p_i^T = frequency of reference allele at locus i , (ancestral) population T .

Under HWE:

$$\Pr(x_{ij} = 2 | p_i^T) = (p_i^T)^2,$$

$$\Pr(x_{ij} = 1 | p_i^T) = 2p_i^T (1 - p_i^T),$$

$$\Pr(x_{ij} = 0 | p_i^T) = (1 - p_i^T)^2.$$

Hardy-Weinberg Equilibrium (HWE): Binomial draws

x_{ij} = genotype at locus i for individual j .

p_i^T = frequency of reference allele at locus i , (ancestral) population T .

Under HWE:

$$\Pr(x_{ij} = 2 | p_i^T) = (p_i^T)^2,$$

$$\Pr(x_{ij} = 1 | p_i^T) = 2p_i^T (1 - p_i^T),$$

$$\Pr(x_{ij} = 0 | p_i^T) = (1 - p_i^T)^2.$$

HWE not valid under population structure!

Goal: measure dependence structure of genotype matrix columns

	Individuals						
Loci	0	2	2	1	1	0	1
	0	2	1	0	1		
	2	...					

X

High-dimensional binomial data

Goal: measure dependence structure of genotype matrix columns

	Individuals						
Loci	0	2	2	1	1	0	1
	0	2	1	0	1		
	2	...					

X

High-dimensional binomial data

Population structure

⇒ dependence between individuals (columns)

Goal: measure dependence structure of genotype matrix columns

	Individuals						
Loci	0	2	2	1	1	0	1
	0	2	1	0	1		
	2	...					

X

High-dimensional binomial data

Population structure

⇒ dependence between individuals (columns)

Linkage disequilibrium

⇒ dependence between loci (rows)

Model parameters

IBD(T): “Identical By Descent” for ancestral population T — shared coin flips

Model parameters

IBD(T): “Identical By Descent” for ancestral population T — shared coin flips

f_j^T : **Inbreeding coefficient**

Pr. that the two alleles at a random locus of individual j are IBD(T)

$$\text{Var}(x_{ij} | T) = 2p_i^T (1 - p_i^T) (1 + f_j^T)$$

Model parameters

IBD(T): “Identical By Descent” for ancestral population T — shared coin flips

f_j^T : **Inbreeding coefficient**

Pr. that the two alleles at a random locus of individual j are IBD(T)

$$\text{Var}(x_{ij} | T) = 2p_i^T (1 - p_i^T) (1 + f_j^T)$$

φ_{jk}^T : **Kinship coefficient**

Pr. that two alleles, one at random from each of individuals j and k , at one random locus are IBD(T)

$$\text{Cov}(x_{ij}, x_{ik} | T) = 4p_i^T (1 - p_i^T) \varphi_{jk}^T$$

Model parameters

IBD(T): “Identical By Descent” for ancestral population T — shared coin flips

f_j^T : **Inbreeding coefficient**

Pr. that the two alleles at a random locus of individual j are IBD(T)

$$\text{Var}(x_{ij} | T) = 2p_i^T (1 - p_i^T) (1 + f_j^T)$$

φ_{jk}^T : **Kinship coefficient**

Pr. that two alleles, one at random from each of individuals j and k , at one random locus are IBD(T)

$$\text{Cov}(x_{ij}, x_{ik} | T) = 4p_i^T (1 - p_i^T) \varphi_{jk}^T$$

F_{ST} : **Fixation index**

Pr. that two random alleles in a subpopulation at a random locus are IBD(T)

Existing approaches

1. F_{ST} estimation

- ▶ *For independent subpopulations only!*
- ▶ Weir-Cockerham (WC) estimator (1984) — 15K citations!
- ▶ “Hudson” pairwise estimator (2013) tweaks WC
- ▶ BayeScan (2008) — 1.2K citations

Existing approaches

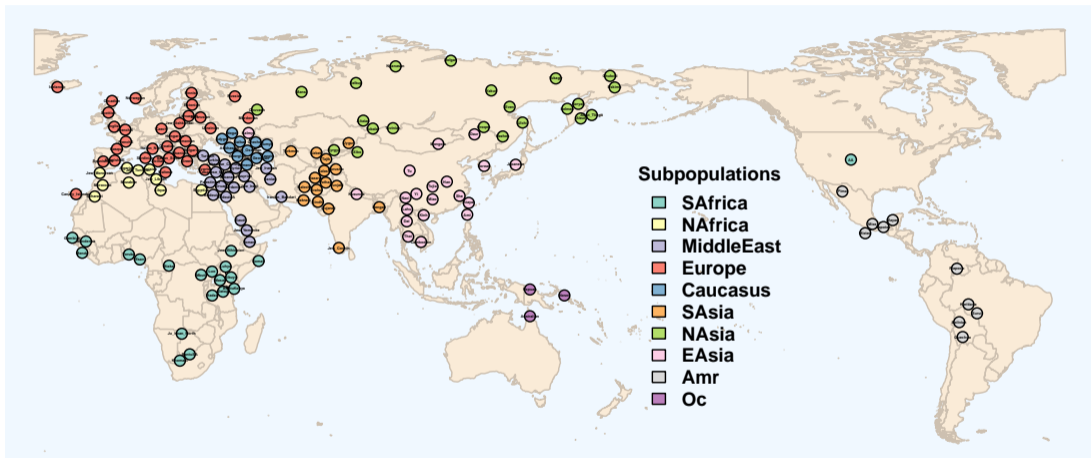
1. F_{ST} estimation

- ▶ *For independent subpopulations only!*
- ▶ Weir-Cockerham (WC) estimator (1984) — 15K citations!
- ▶ “Hudson” pairwise estimator (2013) tweaks WC
- ▶ BayeScan (2008) — 1.2K citations

2. Kinship estimation

- ▶ “Standard” kinship estimator (1950s)
 - ▶ Used by most modern GWAS approaches that control for population structure (PCA, LMM, adj. χ^2 ; top paper 6K citations)
 - ▶ GCTA heritability estimation (2 papers: 4K citations)
- ▶ Our novel finding: accuracy requires unstructured population (a minority of closely-related individuals)

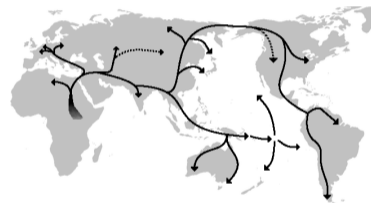
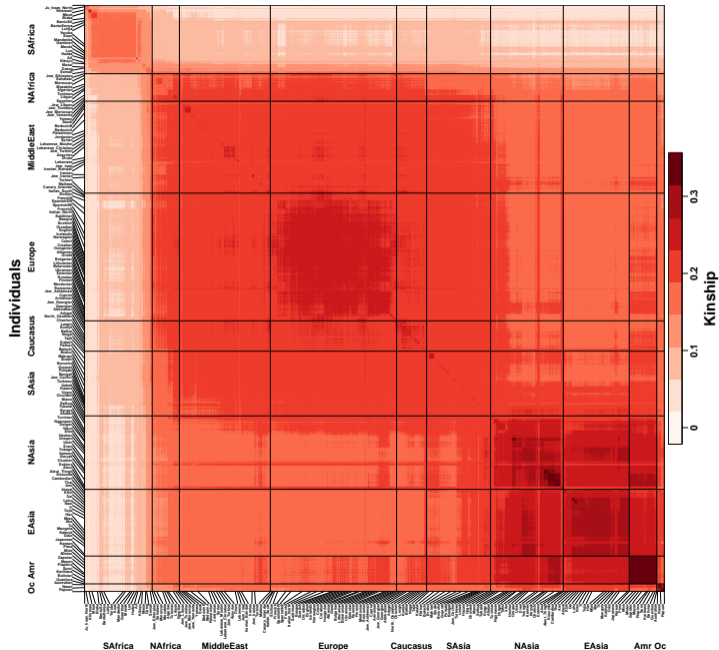
Dataset: Human Origins (Lazaridis *et al.* 2014, 2016)



2,066 indivs. from 163 locs. — 595,911 loci — SNP chip

Our new kinship estimates

Genotypes from "Human Origins"
(Lazaridis *et al.* 2014, 2016)

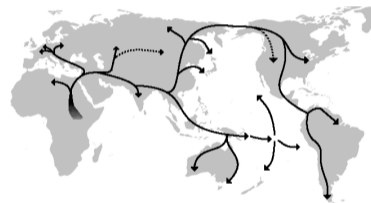
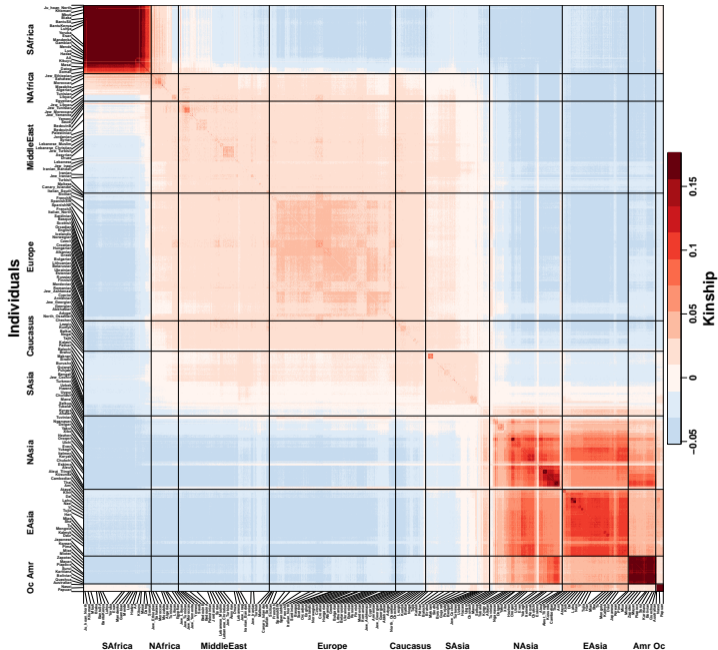


Edited from Ephert [CC BY-SA 3.0], via
Wikimedia Commons

*Inbreeding coeffs. on diagonal

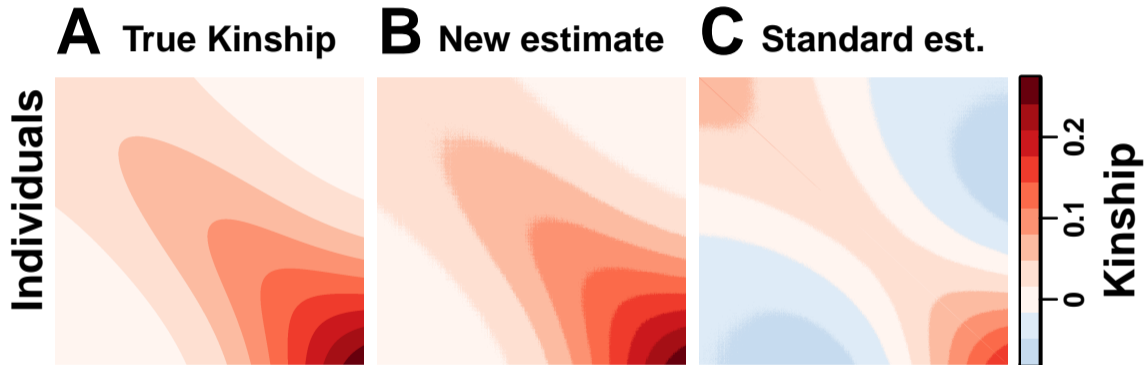
Standard kinship estimates

Genotypes from "Human Origins"
(Lazaridis *et al.* 2014, 2016)

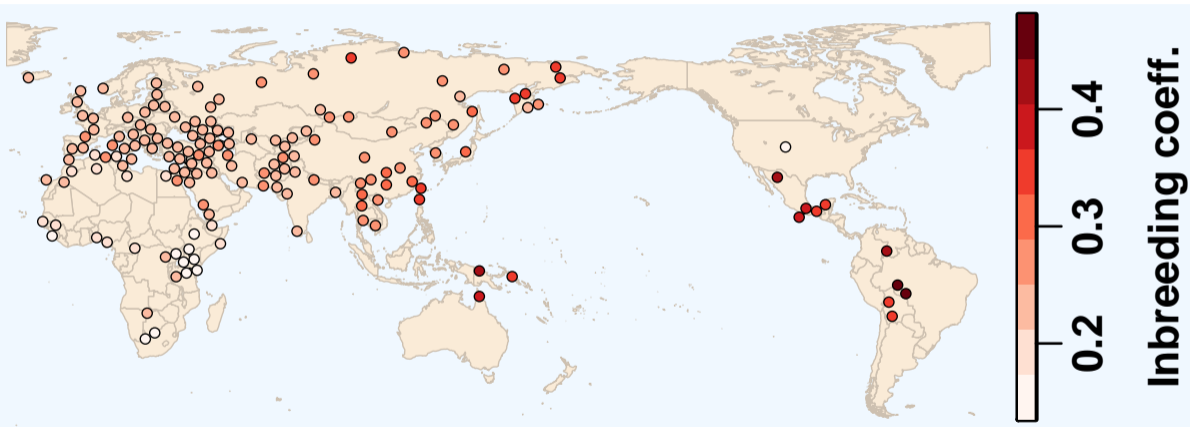


Edited from Ephert [CC BY-SA 3.0], via
Wikimedia Commons

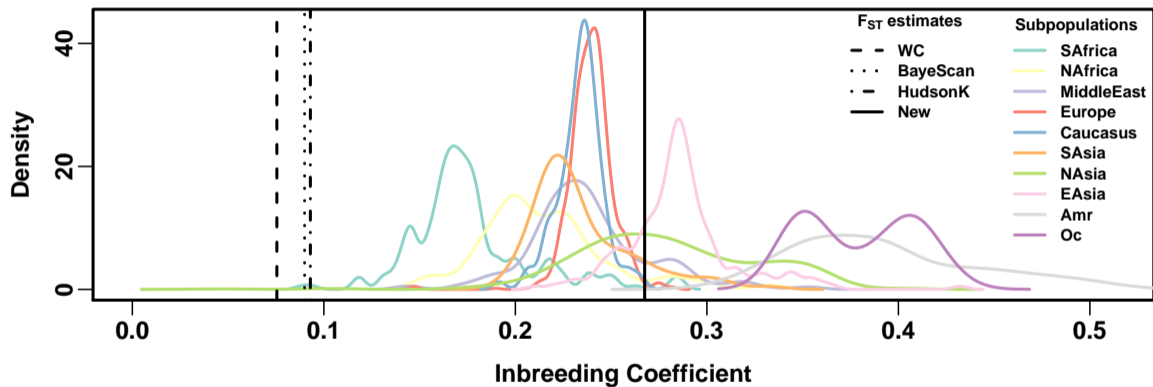
Only our new estimator is accurate in simulations



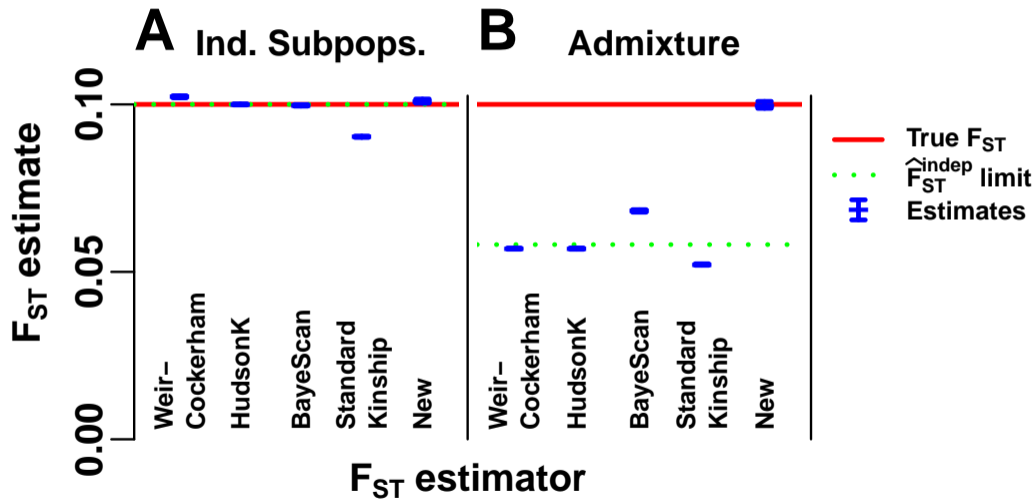
Population-level inbreeding increases with distance from Africa



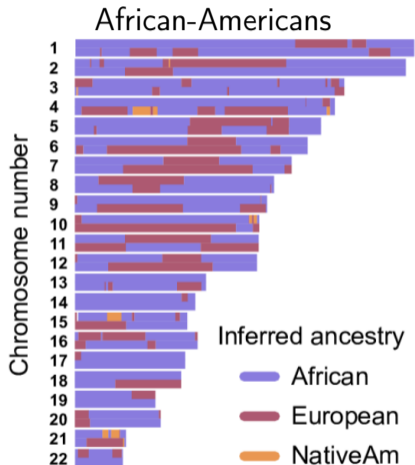
Differentiation (F_{ST}) previously underestimated



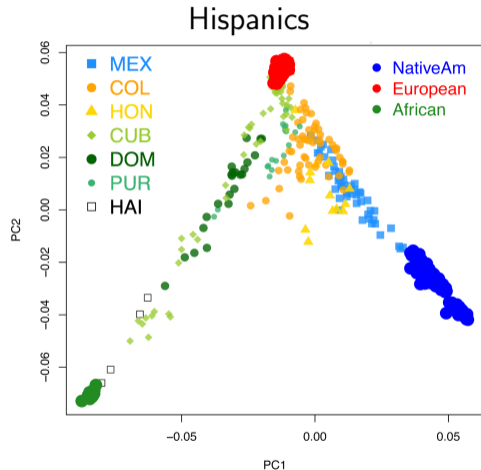
Only our new method estimates generalized F_{ST} accurately



Recently-admixed populations



Baharian *et al.* (2016)



Moreno-Estrada *et al.* (2013)

Admixed siblings from different populations?



Lucy and Maria, UK

Admixed siblings from different populations?



Lucy and Maria, UK



Ochoa brothers, MX

Admixed siblings from different populations?

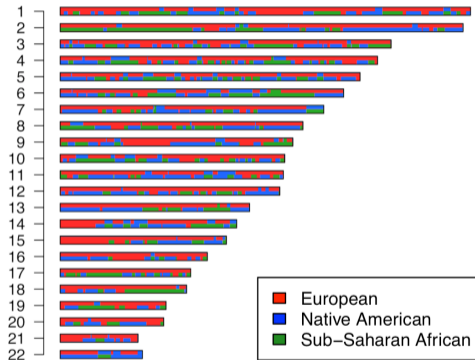


Lucy and Maria, UK



Ochoa brothers, MX

High Admixture LD:



Moreno-Estrada *et al.* (2013)

Admixed siblings from different populations?



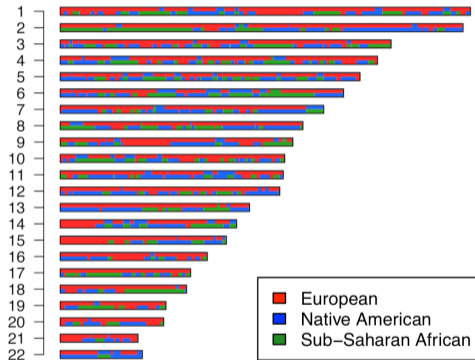
Lucy and Maria, UK



Ochoa brothers, MX

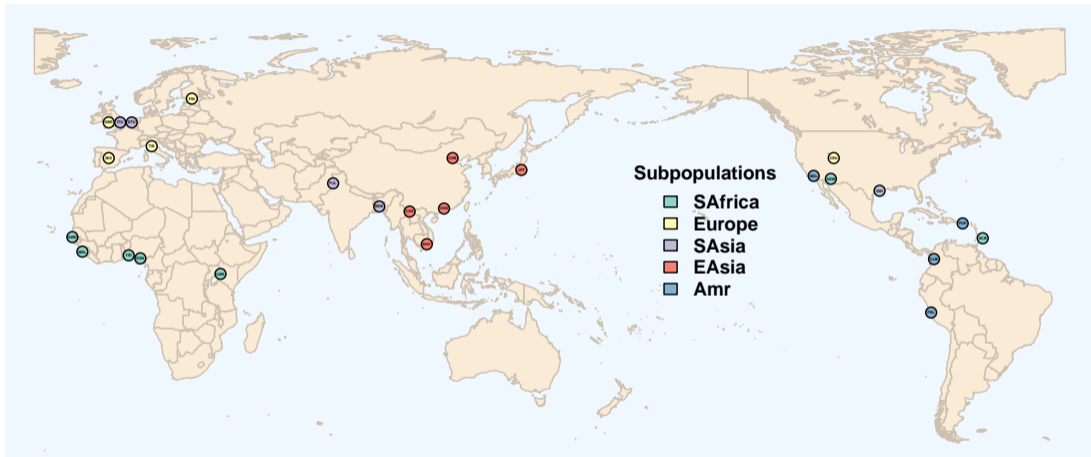
Solution: treat every individual as its own population!

High Admixture LD:



Moreno-Estrada *et al.* (2013)

Dataset: 1000 Genomes Project (2013)

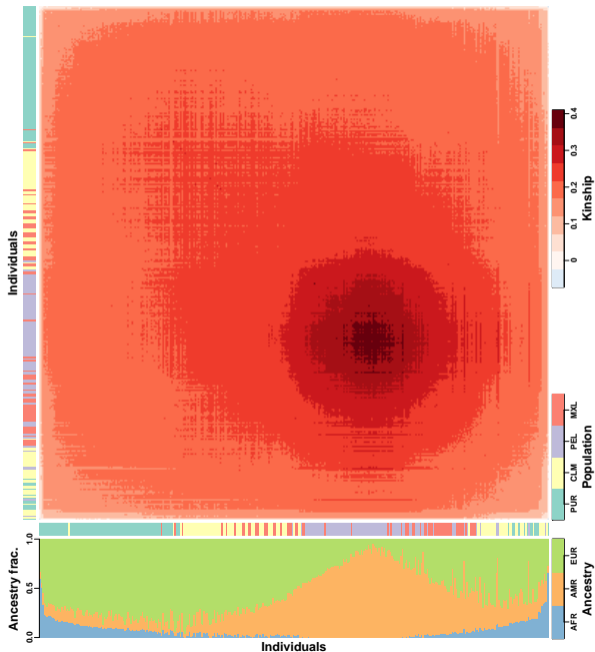


2,504 indivs. from 26 locs. — 20,417,698 loci (asc. in YRI) — WGS trios, etc.

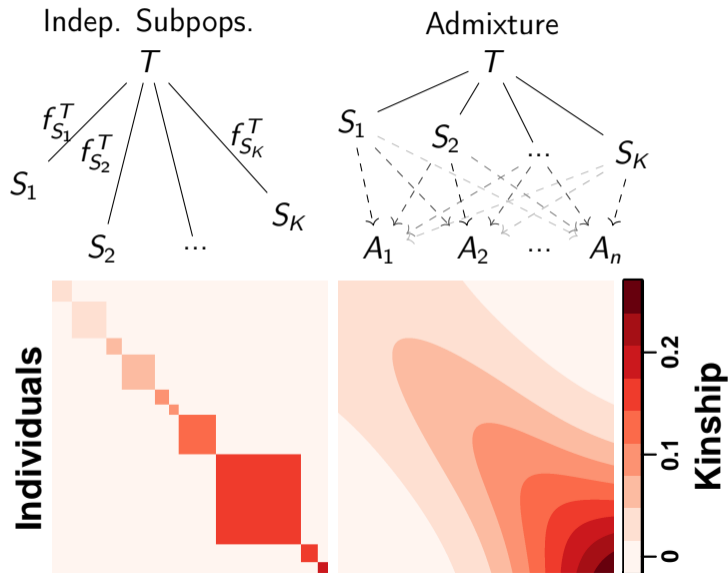
Kinship driven by admixture in Hispanics

Our new kinship estimates

Genotypes from the 1000 Genomes Project (2013)



Comparison of population structures in simulation



F_{ST} in the independent subpopulation model

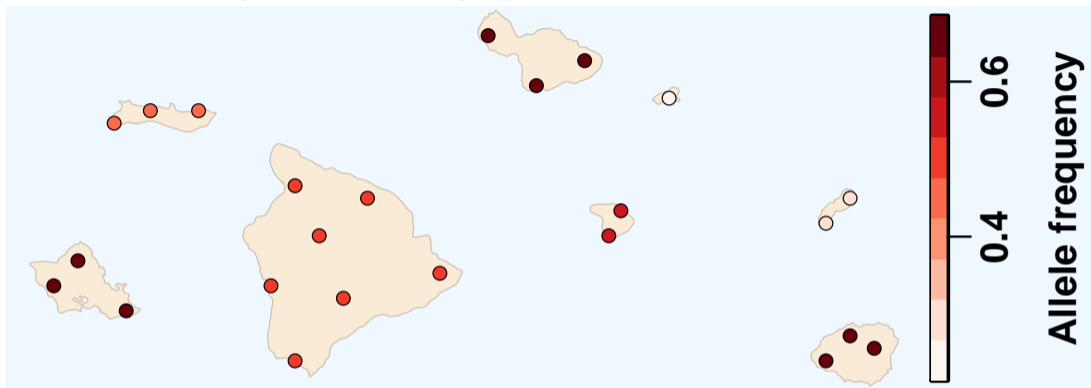


Illustration.

F_{ST} in the independent subpopulation model

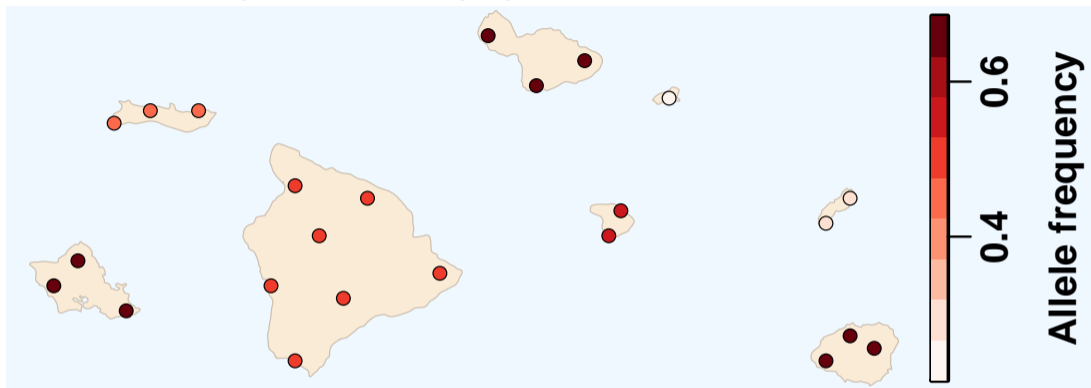


Illustration.

$$F_{ST} = \frac{\text{Var}(p_i^S | T)}{p_i^T (1 - p_i^T)}.$$

Here F_{ST} = proportion of variance explained by pop. structure

Wright's F_{ST}

T = Total, S = Subpopulation, I = Individual.

Total inbreeding: $F_{IT} = \frac{1}{|S|} \sum_{j \in S} f_j^T,$

Local inbreeding: $F_{IS} = \frac{1}{|S|} \sum_{j \in S} f_j^S,$

Structural inbreeding: $F_{ST} = \frac{F_{IT} - F_{IS}}{1 - F_{IS}}.$

Our generalized F_{ST}

Need new “local” subpopulations L_j (separates total from local inbreeding):

$$(1 - f_j^T) = (1 - f_j^{L_j}) (1 - f_{L_j}^T).$$

Our generalized F_{ST}

Need new “local” subpopulations L_j (separates total from local inbreeding):

$$(1 - f_j^T) = (1 - f_j^{L_j}) (1 - f_{L_j}^T).$$

Generalized F_{ST} : applicable to arbitrary population structures, equals previous definition for non-overlapping subpopulations:

$$F_{ST} = \sum_{j=1}^n w_j f_{L_j}^T.$$

Our generalized F_{ST}

Need new “local” subpopulations L_j (separates total from local inbreeding):

$$(1 - f_j^T) = (1 - f_j^{L_j}) (1 - f_{L_j}^T).$$

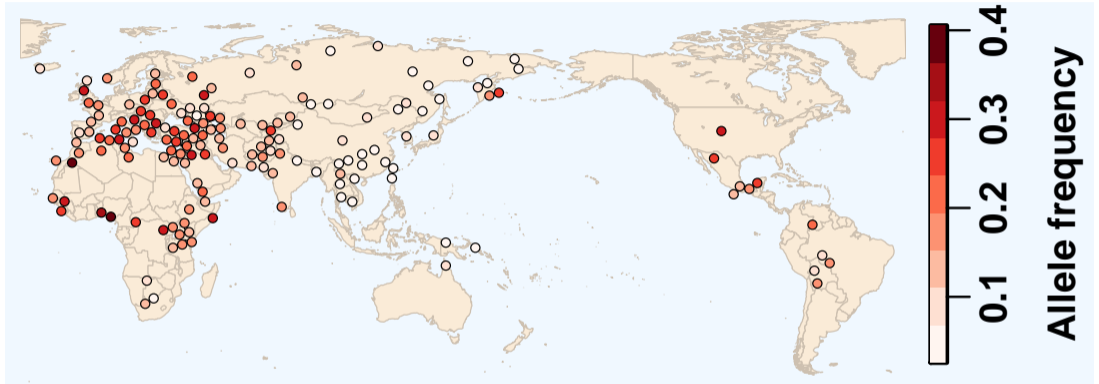
Generalized F_{ST} : applicable to arbitrary population structures, equals previous definition for non-overlapping subpopulations:

$$F_{ST} = \sum_{j=1}^n w_j f_{L_j}^T.$$

Mean heterozygosity in a structured population:

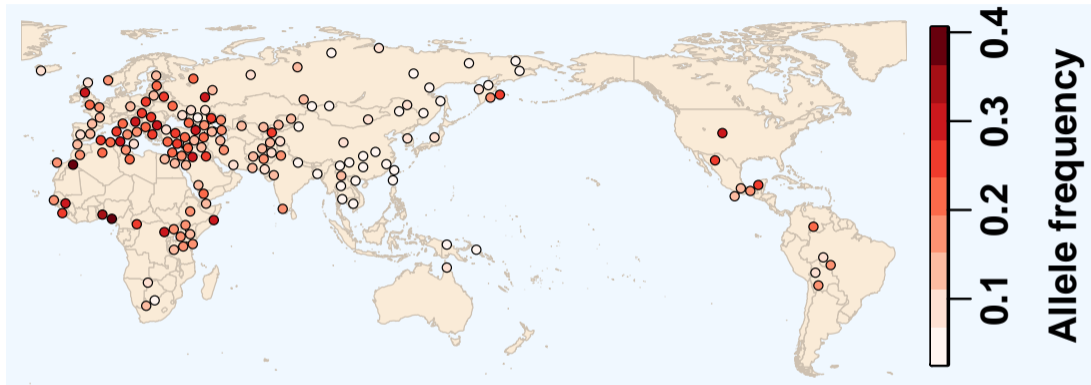
$$\bar{H}_i = \frac{1}{n} \sum_{j=1}^n \Pr(x_{ij} = 1 | T) = 2p_i^T (1 - p_i^T) (1 - F_{ST}).$$

F_{ST} measures population structure / differentiation



Median diff. SNP in Human Origins (rs7626601; given MAF $\geq 10\%$).

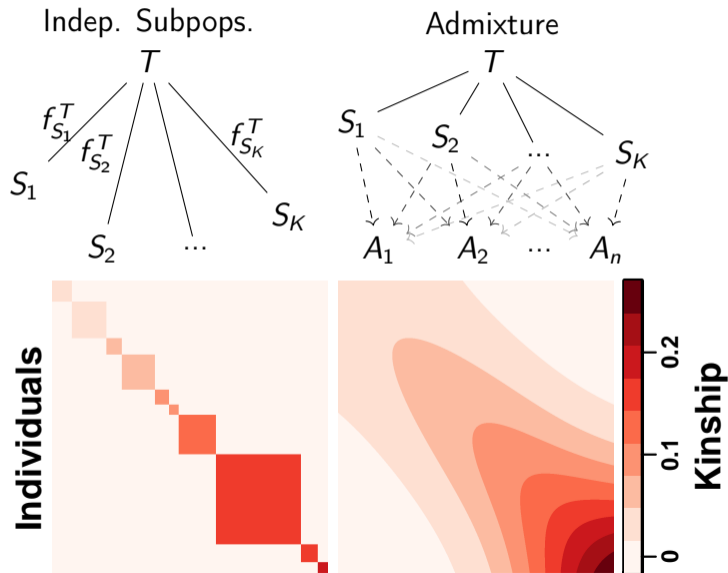
F_{ST} measures population structure / differentiation



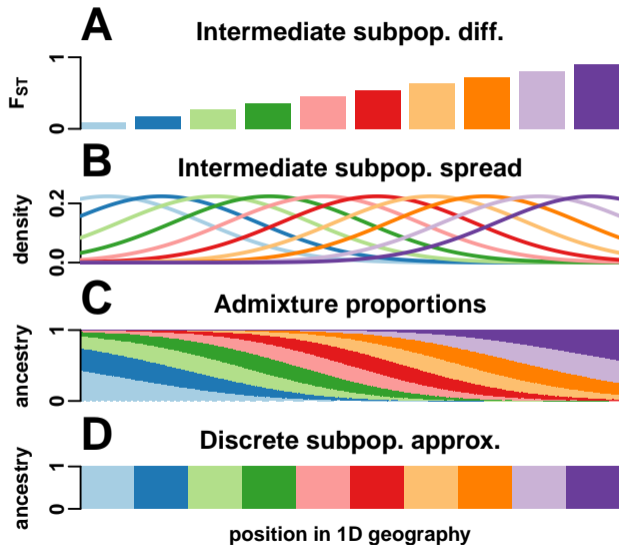
Median diff. SNP in Human Origins (rs7626601; given MAF $\geq 10\%$).

$\hat{F}_{ST}^{WC} \approx 0.0712$ using Weir-Cockerham estimator and $K = 163$.

Comparison of population structures in simulation



Our admixture simulation (R package 'bnpsd' on CRAN)



Kinship model for genotypes

symbol	meaning
T	ref ancestral population
i	locus index
j, k	individual indexes
p_i^T	ref allele frequency
x_{ij}	genotype (num ref alleles)
φ_{jk}^T	kinship of j, k
f_j^T	inbreeding of j

Statistical model:

$$E[x_{ij}|T] = 2p_i^T,$$

$$\text{Var}(x_{ij}|T) = 2p_i^T (1 - p_i^T) (1 + f_j^T),$$

$$\text{Cov}(x_{ij}, x_{ik}|T) = 4p_i^T (1 - p_i^T) \varphi_{jk}^T.$$

(Wright 1921, 1951; Malécot 1948; Jacquard 1970).

Problem: common estimators not consistent under structure

Estimate of ancestral allele frequency:

$$\hat{p}_i^T = \frac{1}{2} \sum_{j=1}^n w_j x_{ij}$$

Variance asymptotically non-zero under population structure:

$$\text{Var}(\hat{p}_i^T | T) = p_i^T (1 - p_i^T) \bar{\varphi}^T$$

Therefore, naive estimators that use \hat{p}_i^T (next) are not consistent!

Bias in standard kinship estimator

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

Bias in standard kinship estimator

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

Bias varies by j, k :

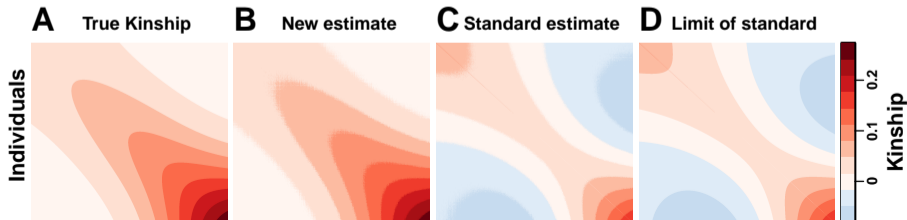
$$\hat{\varphi}_{jk}^{T,\text{std}} \xrightarrow[m \rightarrow \infty]{\text{a.s.}} \frac{\varphi_{jk}^T - \bar{\varphi}_j^T - \bar{\varphi}_k^T + \bar{\varphi}^T}{1 - \bar{\varphi}^T}.$$

Bias in standard kinship estimator

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

Bias varies by j, k :

$$\hat{\varphi}_{jk}^{T,\text{std}} \xrightarrow[m \rightarrow \infty]{\text{a.s.}} \frac{\varphi_{jk}^T - \bar{\varphi}_j^T - \bar{\varphi}_k^T + \bar{\varphi}^T}{1 - \bar{\varphi}^T}.$$



Our new estimator (R package 'popkin' on CRAN)

Step 1: “pre-adjusted” kinship estimator with uniform bias.

$$\hat{\varphi}_{jk}^{T, \text{preadj}} = \frac{\sum_{i=1}^m (x_{ij} - 1)(x_{ik} - 1) - 1}{4 \sum_{i=1}^m \hat{\rho}_i^T (1 - \hat{\rho}_i^T)} + 1 \xrightarrow[m \rightarrow \infty]{\text{a.s.}} \frac{\varphi_{jk}^T - \bar{\varphi}^T}{1 - \bar{\varphi}^T},$$

Our new estimator (R package 'popkin' on CRAN)

Step 1: “pre-adjusted” kinship estimator with uniform bias.

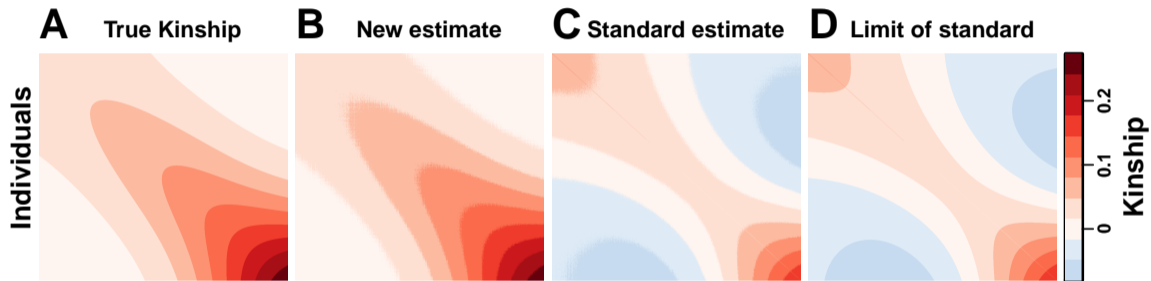
$$\hat{\varphi}_{jk}^{T,\text{preadj}} = \frac{\sum_{i=1}^m (x_{ij} - 1)(x_{ik} - 1) - 1}{4 \sum_{i=1}^m \hat{\rho}_i^T (1 - \hat{\rho}_i^T)} + 1 \xrightarrow[m \rightarrow \infty]{\text{a.s.}} \frac{\varphi_{jk}^T - \bar{\varphi}^T}{1 - \bar{\varphi}^T},$$

Step 2: Estimate minimum kinship, use to unbias “step 1” estimates.

$$\hat{\varphi}_{\min}^{T,\text{preadj}} \xrightarrow[m \rightarrow \infty]{\text{a.s.}} -\frac{\bar{\varphi}^T}{1 - \bar{\varphi}^T}, \quad \hat{\varphi}_{jk}^{T,\text{new}} = \frac{\hat{\varphi}_{jk}^{T,\text{preadj}} - \hat{\varphi}_{\min}^{T,\text{preadj}}}{1 - \hat{\varphi}_{\min}^{T,\text{preadj}}} \xrightarrow[m \rightarrow \infty]{\text{a.s.}} \varphi_{jk}^T.$$

This yields consistent $\hat{f}_j^{T,\text{new}}$, $\hat{F}_{ST}^{\text{new}}$ estimators!

Performance of new estimator



Bias in F_{ST} estimators for independent subpopulations

Previous estimator for n subpopulations, simplified for known AFs (π_{ij}):

$$\hat{F}_{ST}^{\text{indep}} = \frac{\sum_{i=1}^m \hat{\sigma}_i^2}{\sum_{i=1}^m \hat{p}_i^T (1 - \hat{p}_i^T) + \frac{1}{n} \hat{\sigma}_i^2},$$

$$\hat{p}_i^T = \frac{1}{n} \sum_{j=1}^n \pi_{ij}, \quad \hat{\sigma}_i^2 = \frac{1}{n-1} \sum_{j=1}^n (\pi_{ij} - \hat{p}_i^T)^2.$$

Bias in F_{ST} estimators for independent subpopulations

Previous estimator for n subpopulations, simplified for known AFs (π_{ij}):

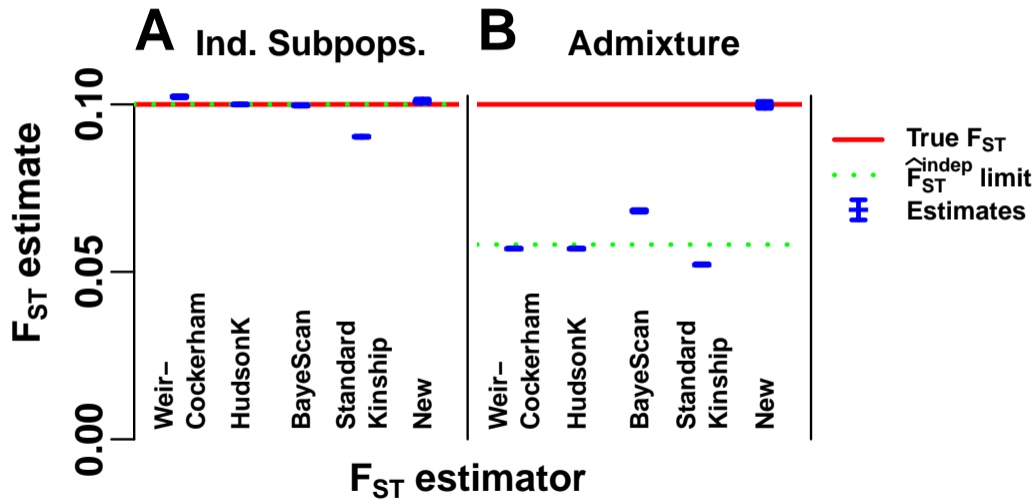
$$\hat{F}_{ST}^{\text{indep}} = \frac{\sum_{i=1}^m \hat{\sigma}_i^2}{\sum_{i=1}^m \hat{p}_i^T (1 - \hat{p}_i^T) + \frac{1}{n} \hat{\sigma}_i^2},$$

$$\hat{p}_i^T = \frac{1}{n} \sum_{j=1}^n \pi_{ij}, \quad \hat{\sigma}_i^2 = \frac{1}{n-1} \sum_{j=1}^n (\pi_{ij} - \hat{p}_i^T)^2.$$

Estimator is biased in dependent subpopulations:

$$\hat{F}_{ST}^{\text{indep}} \xrightarrow[m \rightarrow \infty]{\text{a.s.}} \frac{F_{ST} - \frac{1}{n-1} (n\bar{\theta}^T - F_{ST})}{1 - \frac{1}{n-1} (n\bar{\theta}^T - F_{ST})}.$$

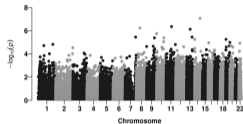
Only our new method estimates generalized F_{ST} accurately



The future: improved kinship has repercussions across genetics!



Accurate and efficient estimation, admixture modeling



Association studies, selection tests



Bias in heritability of complex traits



Animal and plant breeding

Acknowledgments

John D. Storey

Andrew Bass

Irineo Cabrerros

Wei Hao

Riley Skeen-Gaar

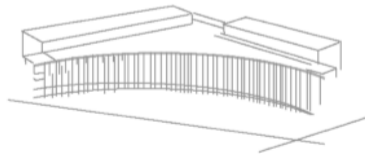
Neo Christopher Chung

University of Warsaw

Funding:

National Institutes of Health

Otsuka Pharmaceutical



Lewis-Sigler Institute for Integrative Genomics