

# Challenges and new approaches for whole genome analysis in multi-ethnic studies with pedigrees.

Timothy Thornton, PhD

Robert W. Day Endowed Professor of Public Health

Department of Biostatistics

University of Washington

## New Statistical Methods for Family-Based Sequencing Studies

BIRS: August 5-10, 2018



SCHOOL OF PUBLIC HEALTH  
UNIVERSITY of WASHINGTON

# Motivation: Multi-Ethnic PAGE Study

- The Population Architecture using Genomics and Epidemiology (PAGE) study focuses on exploring the genetics of underrepresented populations.
- PAGE consists of **49,839** individuals of non-European ancestry from multi-ethnic studies,
- Individuals were genotyped using the approximately 1.3 million variants on the Multi-Ethnic Global Array (MEGA) imputed to 1000 Genomes Phase 3

# Multi-Ethnic PAGE Study

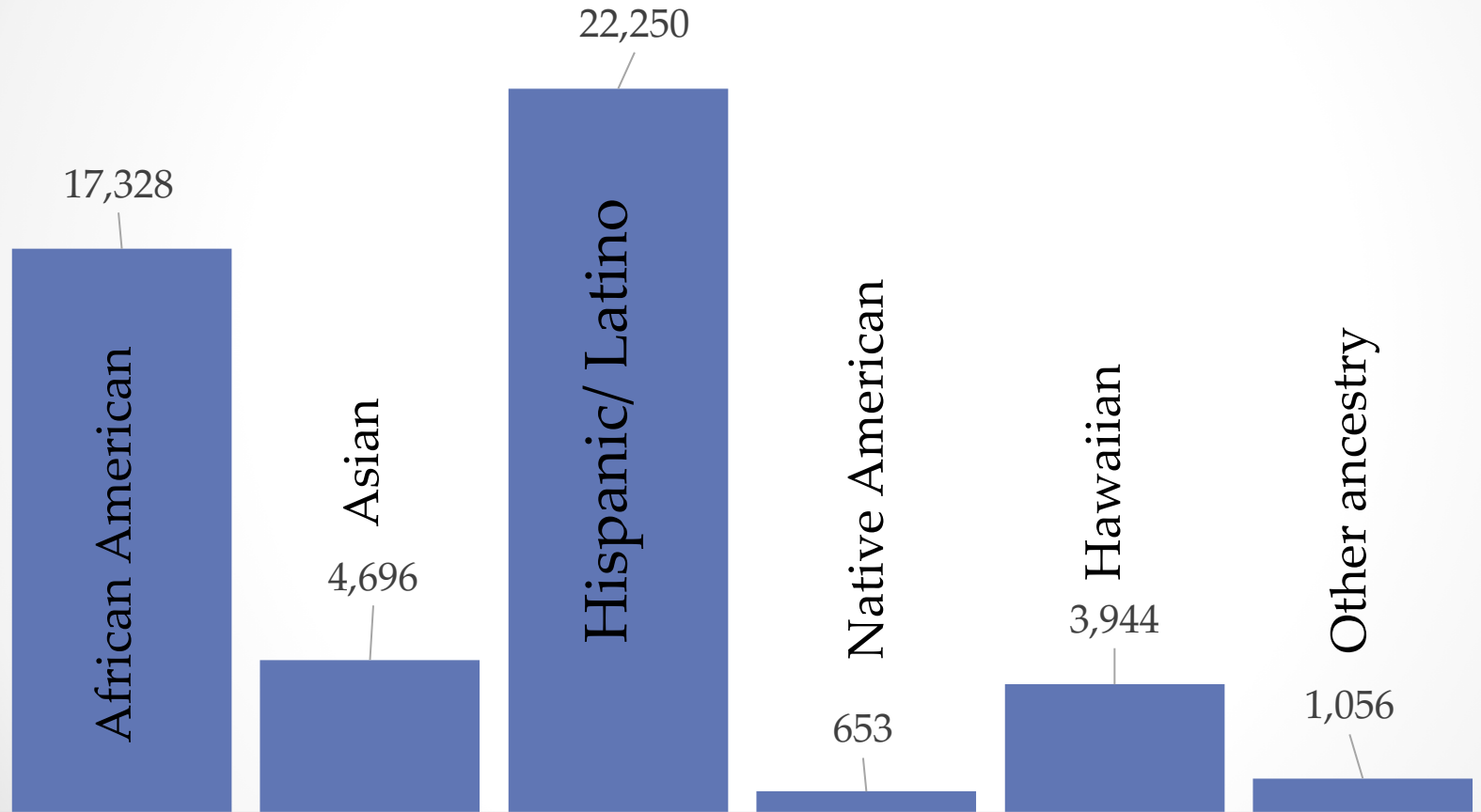


Figure Courtesy of TOPMed DCC

# Multi-Ethnic PAGE Study

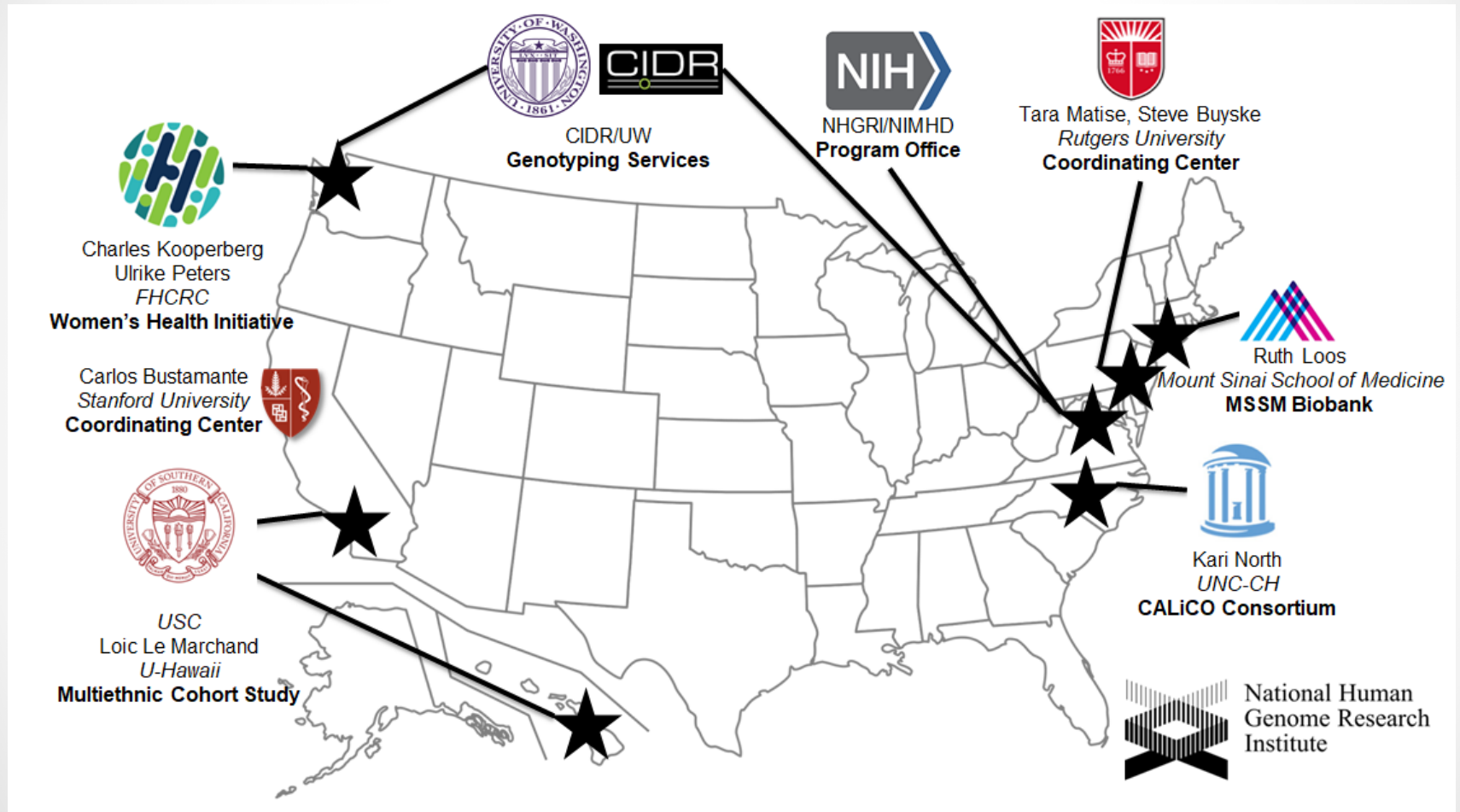


Figure Courtesy of Mariaelisa Graff (UNC Chapel Hill)

# Motivation: TOPMed WGS Project

- NIH/NHLBI Trans-Omics for Precision Medicine (TOPMed) Whole-genome-sequence (WGS) project is generated deep whole genome sequencing data for more than **170,000 individuals**.
- Cohorts are from multi-ethnic populations with well-defined phenotypes and existing clinical outcomes data.

# TOPMed WGS Project: Multi-Ethnic Cohorts

- 30X coverage Illumina X-10 sequencing
- 4 sequencing centers with harmonized protocols

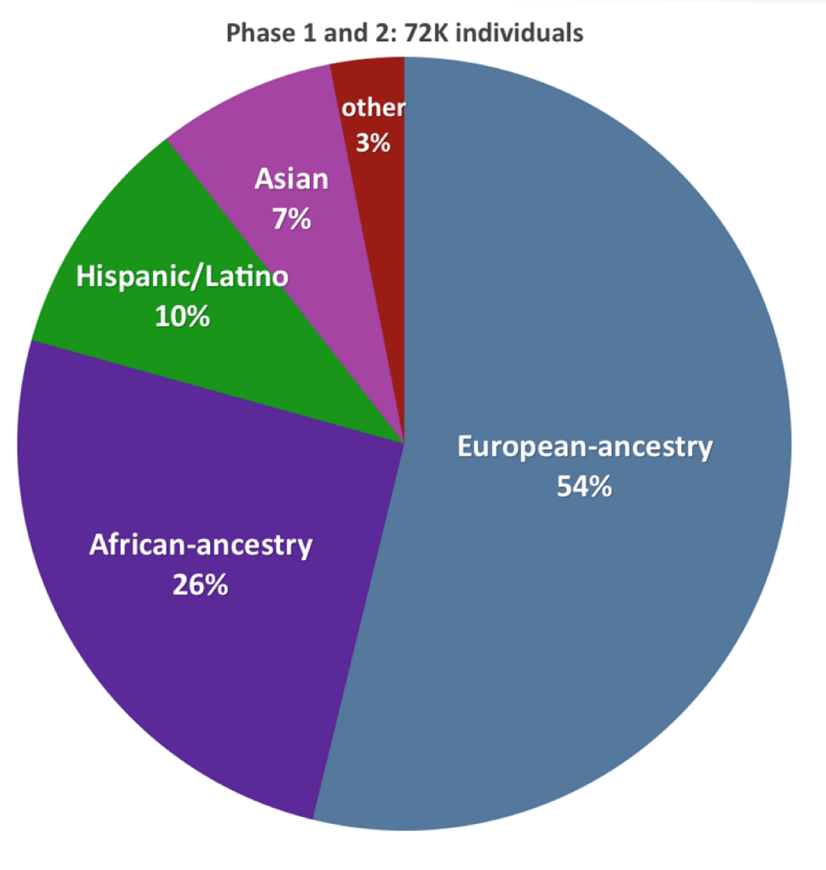


Chart from TOPMed DCC

# PAGE and TOPMed and PAGE: Opportunities

- The PAGE and TOPMed WGS Project offer unprecedented opportunities for:
  - Identification of population-specific variants as well as novel low frequency and rare genetic variants underlying phenotypic diversity
  - Potential to provide new insights into human health and health disparities of minority populations for many complex diseases



# TOPMed and PAGE : Challenges

- Challenges for analysis of PAGE and TOPMed whole genome data:
  - **Multi-ethnic populations and confounding due to highly heterogeneous genetic and environmental backgrounds within and across cohorts**
  - **Variety of study designs: case-control and cohort studies, family-based studies, founder populations (e.g., TOPMed includes the Framingham Heart Study, Jackson Heart Study, and Amish samples).**
  - Computational burden for analysis of deep whole genome sequence data for 120,000+ individuals (TOPMed) and 50,000 individuals with 1000 Genomes imputed genotypes (PAGE)



# Linear Mixed Models

- Linear mixed models (LMMs) have emerged as a powerful and effective method of choice for genetic association testing in the presence of sample structure
- LMMs have been used to simultaneously account for both population structure, family structure, and/or cryptic relatedness

# Linear Mixed Models

- A number of LMMs have been proposed including EMMAX proposed by Kang et al. [Nat Genet, 2010], GEMMA proposed by Zhou and Stephens [Nat Genet, 2012] and others:

## TECHNICAL REPORTS

nature  
genetics

Variance component model to account for sample structure in genome-wide association studies

Hyun Min Kang<sup>1,2,8</sup>, Jae Hoon Sul<sup>3,8</sup>, Susan K Service<sup>4</sup>, Noah A Zaitlen<sup>5</sup>, Sit-yeec Kong<sup>4</sup>, Nelson B Freimer<sup>4</sup>, Chiara Sabatti<sup>6</sup> & Eleazar Eskin<sup>3,7</sup>

## TECHNICAL REPORTS

nature  
genetics

Rapid variance components-based method for whole-genome association analysis

Gulnara R Svishcheva<sup>1</sup>, Tatiana I Axenovich<sup>1</sup>, Nadezhda M Belonogova<sup>1</sup>, Cornelia M van Duijn<sup>2</sup> & Yuri S Aulchenko<sup>1</sup>

## TECHNICAL REPORTS

nature  
genetics

Genome-wide efficient mixed-model analysis for association studies

Xiang Zhou<sup>1</sup> & Matthew Stephens<sup>1,2</sup>

## TECHNICAL REPORTS

nature  
genetics

Mixed linear model approach adapted for genome-wide association studies

Zhiwu Zhang<sup>1</sup>, Elhan Ersoz<sup>1</sup>, Chao-Qiang Lai<sup>2</sup>, Rory J Todhunter<sup>3</sup>, Hemant K Tiwari<sup>4</sup>, Michael A Gore<sup>5</sup>, Peter J Bradbury<sup>6</sup>, Jianming Yu<sup>7</sup>, Donna K Arnett<sup>8</sup>, Jose M Ordovas<sup>2,9</sup> & Edward S Buckler<sup>1,6</sup>

# LMMs for Genetic Association Testing

- Most LMM methods fit the following linear mixed model

$$\mathbf{Y} = \mathbf{g}_s \beta_s + \mathbf{X} \boldsymbol{\alpha} + \boldsymbol{\epsilon} \quad \text{with} \quad \boldsymbol{\epsilon} \sim N(\mathbf{0}, \boldsymbol{\Sigma} \equiv \sigma_A^2 \boldsymbol{\Psi} + \sigma_\epsilon^2 \mathbf{I})$$

where:

- $\mathbf{Y}$  is the vector of phenotype values
- $\mathbf{g}_s$  is the vector of genotypes at the SNP being tested,
- $\beta_s$  is the (scalar) association parameter of interest, measuring the effect of genotype on phenotype
- $\mathbf{X}$  is a matrix of covariate values with vector  $\boldsymbol{\alpha}$  of covariate effects
- $\boldsymbol{\Psi}$  is an estimated genetic relationship matrix (GRM) capturing population structure and relatedness
- $\sigma_A^2$  is the additive genetic variance for polygenic effects
- $\mathbf{I}$  is the identity matrix
- $\sigma_\epsilon^2$  presents non-genetic variance due to non-genetic effects assumed to be acting independently on individuals

# Standard GRM for LMMs

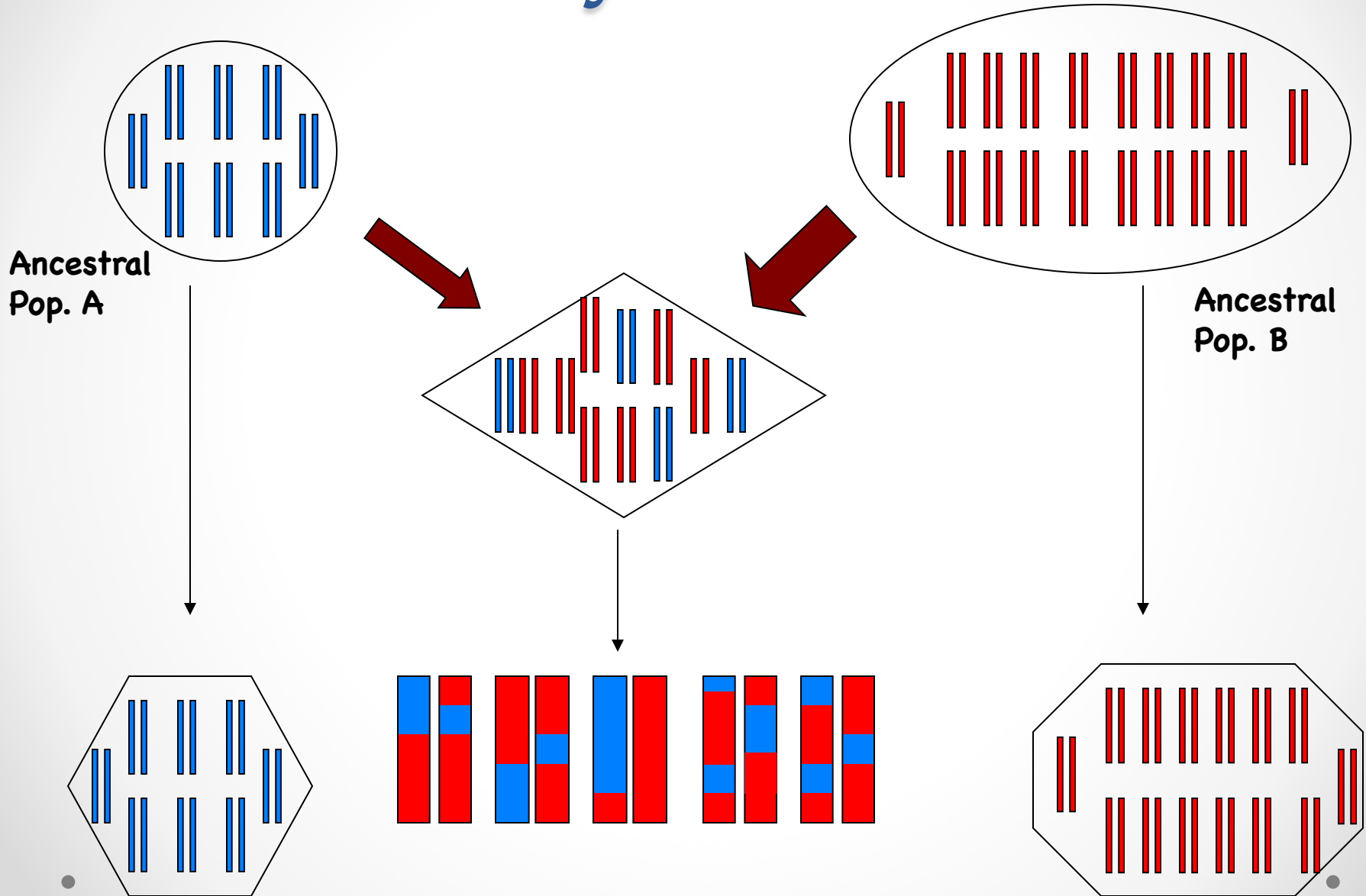
- Average genotypic correlations across the genome is summarized with a single genetic relatedness matrix (GRM)
- GRM entry for subjects  $i$  and  $j$  who have  $M$  genotyped markers across the autosomes:

$$\hat{\Psi}_{ij} = \frac{1}{M} \sum_{m=1}^M \frac{(g_{mi} - 2\hat{p}_m)(g_{mj} - 2\hat{p}_m)}{2\hat{p}_m(1 - \hat{p}_m)}$$

# Linear Mixed Models for Genetic Association Testing

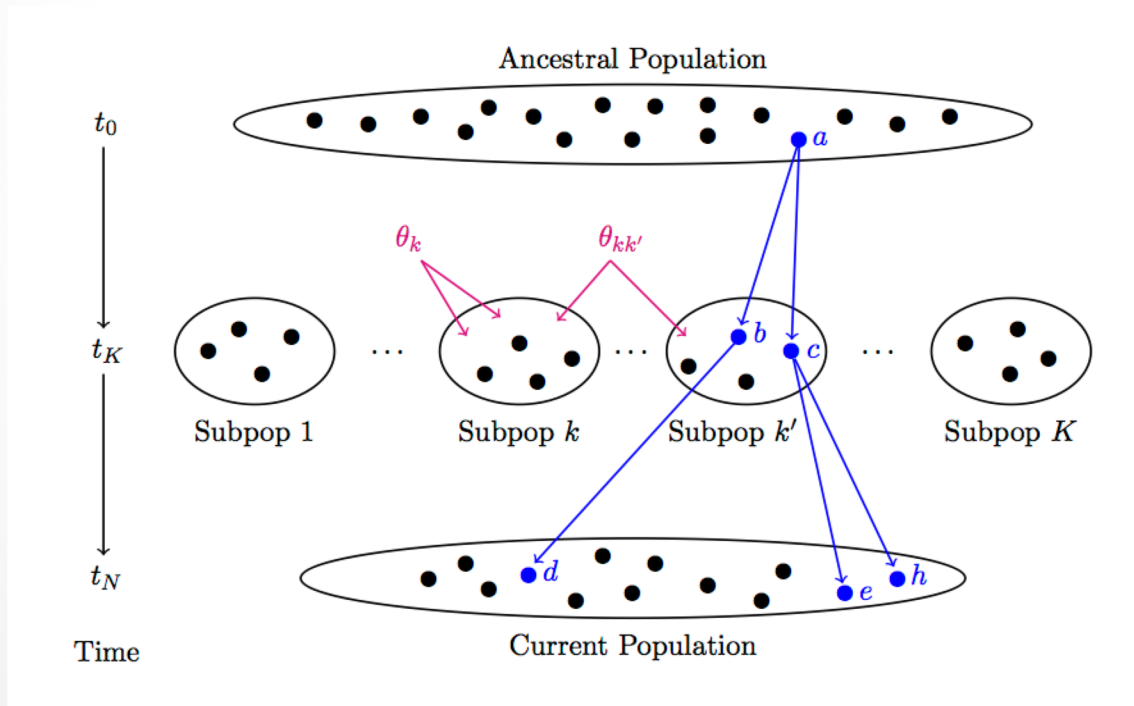
- LMMs have largely been evaluated in samples with relatively subtle population structure, e.g., European populations
- Existing LMMs methods use a single empirical genetic relationship matrix to model the entire genealogy of sampled individuals as part of the covariance structure of the phenotype.
- Samples from multi-ethnic cohorts often have complex genealogy due to ancestry admixture and both recent and distant genetic relatedness
- The genealogy of sampled individuals consists of:
  - Distant genetic relatedness, such as population structure
  - Recent genetic relatedness: pedigree relationships of close relatives

# Ancestry Admixture



# Recent versus Distant Genetic Relatedness

- Distinguishing familial relatedness from ancestry using genotype data in diverse populations is difficult, as both manifest as genetic similarity through the sharing of alleles.



# Deconvolution of Genetic Relatedness

- Conomos et al., *Am J Hum Genet*, 2016

**ARTICLE**

## Model-free Estimation of Recent Genetic Relatedness

Matthew P. Conomos,<sup>1,\*</sup> Alexander P. Reiner,<sup>2,3</sup> Bruce S. Weir,<sup>1</sup> and Timothy A. Thornton<sup>1,\*</sup>

- Conomos et al., *Genet Epidemiology*, 2015

RESEARCH ARTICLE

### **Robust Inference of Population Structure for Ancestry Prediction and Correction of Stratification in the Presence of Relatedness**

Matthew P. Conomos,<sup>1</sup> Michael B. Miller,<sup>2</sup> and Timothy A. Thornton<sup>1\*</sup>

Genetic  
Epidemiology

OFFICIAL JOURNAL  
INTERNATIONAL GENETIC  
EPIDEMIOLOGY SOCIETY  
[www.geneticepi.org](http://www.geneticepi.org)

- Thornton et al., *Am J Hum Genet*, 2012

**ARTICLE**

## Estimating Kinship in Admixed Populations

Timothy Thornton,<sup>1,\*</sup> Hua Tang,<sup>2</sup> Thomas J. Hoffmann,<sup>3,4</sup> Heather M. Ochs-Balcom,<sup>5</sup> Bette J. Caan,<sup>6</sup>  
and Neil Risch<sup>3,4,6,\*</sup>



# LMM-OPS for Multi-Ethnic Populations



- Matt Conomos PhD dissertation work developed LMM-OPS for association mapping in ancestrally diverse populations
- LMM-OPS, linear mixed models with orthogonal partitioned structure
- Appropriately accounts for the complex genealogy of ancestrally diverse samples by partitioning sample structure into two orthogonal components:
  1. a component for the sharing of alleles inherited identical by descent (IBD) from recent common ancestors, which represents familial relatedness
  2. and another component for allele sharing due to more distant common ancestry, which represents population structure.

# LMM-OPS for Multi-Ethnic Populations

- With LMM-OPS, a score test for association is calculated based on the following linear mixed model:

$$\mathbf{Y} = \mathbf{g}_s \beta_s + \mathbf{X} \boldsymbol{\alpha} + \mathbf{V} \boldsymbol{\gamma} + \boldsymbol{\epsilon} \quad \text{with} \quad \boldsymbol{\epsilon} \sim \text{N}(\mathbf{0}, \boldsymbol{\Sigma} \equiv \sigma_A^2 \boldsymbol{\Phi} + \sigma_\epsilon^2 \mathbf{I})$$

where:

- $\boldsymbol{\Phi}$  is an genetic relatedness matrix adjusted for ancestry admixture (via the PCs) with PC-Relate
- $\mathbf{V}$  is a matrix with PCs from PC-AiR, and  $\boldsymbol{\gamma}$  is a vector (unknown) ancestry effects on the phenotype
- $\mathbf{X}$  is a matrix of covariate values with vector  $\boldsymbol{\alpha}$  of covariate effects

# PC-Relate GRMs for LMM-OPS

Two possible GRMs we have considered for LMM-OPS are as follows, with the following entries for subjects  $i$  and  $j$  who have  $M$  genotyped markers.

$$\hat{\Phi}_{ij}^1 = \frac{1}{M} \sum_{m=1}^M \frac{(g_{mi} - 2\hat{p}_{mi})(g_{mj} - 2\hat{p}_{mj})}{\sqrt{2\hat{p}_{mi}(1-\hat{p}_{mi})}\sqrt{2\hat{p}_{mj}(1-\hat{p}_{mj})}}$$

$$\hat{\Phi}_{ij}^2 = \frac{\frac{1}{M} \sum_{m=1}^M (g_{mi} - 2\hat{p}_{mi})(g_{mj} - 2\hat{p}_{mj})}{\frac{1}{M} \sum_{m=1}^M \sqrt{2\hat{p}_{mi}(1-\hat{p}_{mi})}\sqrt{2\hat{p}_{mj}(1-\hat{p}_{mj})}}$$

Both use **individual specific allele** frequencies,  $\hat{p}_{mi}$ , calculated from a regression analysis with the PCs from PC-AiR included as predictors of genotype

# Simulations with Admixture: Genomic Control Inflation

## Evaluation

Method	Genome-Wide	Highly <sup>a</sup> Differentiated	Moderately <sup>b</sup> Differentiated	Weakly <sup>c</sup> Differentiated
<b>LMM-OPS</b>	1.000 (0.0002)	0.999 (0.0007)	1.001 (0.0004)	1.001 (0.0003)
<b>EMMAX</b>	1.001 (0.0002)	<b>1.098 (0.0011)</b>	<b>1.016 (0.0004)</b>	<b>0.979 (0.0003)</b>
<b>GEMMA</b>	1.004 (0.0002)	<b>1.110 (0.0011)</b>	<b>1.020 (0.0005)</b>	<b>0.980 (0.0003)</b>
<b>Linear Reg. with PCs</b>	<b>1.026 (0.0006)</b>	<b>1.025 (0.0009)</b>	<b>1.027 (0.0007)</b>	<b>1.026 (0.0006)</b>

<sup>a</sup> Highly differentiated SNPs:  $D_s \geq 0.4$  between the two populations

<sup>b</sup> Moderately differentiated SNPs:  $0.4 > D_s \geq 0.2$  between the two populations

<sup>c</sup> Weakly differentiated SNPs:  $D_s < 0.2$  between the two populations

# Applications and Discoveries in Hispanic/Latino Populations

## ARTICLE

### Genetic Diversity and Association Studies in US Hispanic/Latino Populations: Applications in the Hispanic Community Health Study/Study of Latinos

Matthew P. Conomos,<sup>1,14,\*</sup> Cecelia A. Laurie,<sup>1,14</sup> Adrienne M. Stilp,<sup>1,14</sup> Stephanie M. Gogarten,<sup>1,14</sup> Caitlin P. McHugh,<sup>1</sup> Sarah C. Nelson,<sup>1</sup> Tamar Sofer,<sup>1</sup> Lindsay Fernández-Rhodes,<sup>2</sup> Anne E. Justice,<sup>2</sup> Mariaelisa Graff,<sup>2</sup> Kristin L. Young,<sup>2</sup> Amanda A. Seyerle,<sup>2</sup> Christy L. Avery,<sup>2</sup> Kent D. Taylor,<sup>3</sup> Jerome I. Rotter,<sup>3</sup> Gregory A. Talavera,<sup>4</sup> Martha L. Daviglius,<sup>5</sup> Sylvia Wassertheil-Smoller,<sup>6</sup> Neil Schneiderman,<sup>7</sup> Gerardo Heiss,<sup>2</sup> Robert C. Kaplan,<sup>6</sup> Nora Franceschini,<sup>2</sup> Alex P. Reiner,<sup>8</sup> John R. Shaffer,<sup>9</sup> R. Graham Barr,<sup>10</sup> Kathleen F. Kerr,<sup>1</sup> Sharon R. Browning,<sup>1</sup> Brian L. Browning,<sup>11</sup> Bruce S. Weir,<sup>1</sup> M. Larissa Avilés-Santa,<sup>12</sup> George J. Papanicolaou,<sup>12</sup> Thomas Lumley,<sup>13</sup> Adam A. Szpiro,<sup>1</sup> Kari E. North,<sup>2</sup> Ken Rice,<sup>1</sup> Timothy A. Thornton,<sup>1</sup> and Cathy C. Laurie<sup>1,\*</sup>

## ARTICLE

### Genome-wide Association Study of Platelet Count Identifies Ancestry-Specific Loci in Hispanic/Latino Americans

Ursula M. Schick,<sup>1,2,3,16</sup> Deepti Jain,<sup>4,16</sup> Chani J. Hodonsky,<sup>5,16</sup> Jean V. Morrison,<sup>4</sup> James P. Davis,<sup>6</sup> Lisa Brown,<sup>4</sup> Tamar Sofer,<sup>4</sup> Matthew P. Conomos,<sup>4</sup> Claudia Schurmann,<sup>2,3</sup> Caitlin P. McHugh,<sup>4</sup> Sarah C. Nelson,<sup>4</sup> Swarooparani Vadlamudi,<sup>6</sup> Adrienne Stilp,<sup>4</sup> Anna Plantinga,<sup>4</sup> Leslie Baier,<sup>7</sup> Stephanie A. Bien,<sup>1</sup> Stephanie M. Gogarten,<sup>4</sup> Cecelia A. Laurie,<sup>4</sup> Kent D. Taylor,<sup>8,9</sup> Yongmei Liu,<sup>10</sup> Paul L. Auer,<sup>11</sup> Nora Franceschini,<sup>5</sup> Adam Szpiro,<sup>4</sup> Ken Rice,<sup>4</sup> Kathleen F. Kerr,<sup>4</sup> Jerome I. Rotter,<sup>8</sup> Robert L. Hanson,<sup>7</sup> George Papanicolaou,<sup>12</sup> Stephen S. Rich,<sup>13,14</sup> Ruth J.F. Loos,<sup>2,3,15</sup> Brian L. Browning,<sup>4</sup> Sharon R. Browning,<sup>4</sup> Bruce S. Weir,<sup>4</sup> Cathy C. Laurie,<sup>4</sup> Karen L. Mohlke,<sup>6</sup> Kari E. North,<sup>5,16</sup> Timothy A. Thornton,<sup>4,16</sup> and Alex P. Reiner<sup>1,16,\*</sup>

## ASSOCIATION STUDIES ARTICLE

### Genome-wide association study of dental caries in the Hispanic Communities Health Study/Study of Latinos (HCHS/SOL)

Jean Morrison<sup>1</sup>, Cathy C. Laurie<sup>1</sup>, Mary L. Marazita<sup>2,3,4</sup>, Anne E. Sanders<sup>5</sup>, Steven Offenbacher<sup>6</sup>, Christian R. Salazar<sup>7,8</sup>, Matthew P. Conomos<sup>1</sup>, Timothy Thornton<sup>1</sup>, Deepti Jain<sup>1</sup>, Cecelia A. Laurie<sup>1</sup>, Kathleen F. Kerr<sup>1</sup>, George Papanicolaou<sup>9</sup>, Kent Taylor<sup>10</sup>, Linda M. Kaste<sup>11</sup>, James D. Beck<sup>5</sup> and John R. Shaffer<sup>2,\*</sup>

## ORIGINAL ARTICLE

### Genetic Associations with Obstructive Sleep Apnea Traits in Hispanic/Latino Americans

Brian E. Cade<sup>1,2</sup>, Han Chen<sup>3</sup>, Adrienne M. Stilp<sup>4</sup>, Kevin J. Gleason<sup>1</sup>, Tamar Sofer<sup>4</sup>, Sonia Ancoli-Israel<sup>6,6,7</sup>, Raanan Arens<sup>8</sup>, Graeme I. Bell<sup>9</sup>, Jennifer E. Below<sup>10</sup>, Andrew C. Bjornes<sup>11</sup>, Sung Chun<sup>11,12</sup>, Matthew P. Conomos<sup>4</sup>, Daniel S. Evans<sup>13</sup>, W. Craig Johnson<sup>4</sup>, Alexis C. Frazier-Wood<sup>14</sup>, Jacqueline M. Lane<sup>1,2,15,16</sup>, Emma K. Larkin<sup>17</sup>, Jose S. Loredó<sup>18</sup>, Wendy S. Post<sup>19</sup>, Alberto R. Ramos<sup>20</sup>, Ken Rice<sup>4</sup>, Jerome I. Rotter<sup>21</sup>, Neomi A. Shah<sup>22</sup>, Katie L. Stone<sup>13</sup>, Kent D. Taylor<sup>21</sup>, Timothy A. Thornton<sup>4</sup>, Gregory J. Tranah<sup>13</sup>, Chaolong Wang<sup>3,23</sup>, Phyllis C. Zee<sup>24</sup>, Craig L. Hanis<sup>10</sup>, Shamil R. Sunyaev<sup>11,12,16</sup>, Sanjay R. Patel<sup>1,2,25</sup>, Cathy C. Laurie<sup>4</sup>, Xiaofeng Zhu<sup>26</sup>, Richa Saxena<sup>1,15,16</sup>, Xihong Lin<sup>3</sup>, and Susan Redline<sup>1,2,25</sup>

## ARTICLE

### Control for Population Structure and Relatedness for Binary Traits in Genetic Association Studies via Logistic Mixed Models

Han Chen,<sup>1,8</sup> Chaolong Wang,<sup>1,2,8</sup> Matthew P. Conomos,<sup>3</sup> Adrienne M. Stilp,<sup>3</sup> Zilin Li,<sup>1,4</sup> Tamar Sofer,<sup>3</sup> Adam A. Szpiro,<sup>3</sup> Wei Chen,<sup>5</sup> John M. Brehm,<sup>5</sup> Juan C. Celedón,<sup>5</sup> Susan Redline,<sup>6</sup> George J. Papanicolaou,<sup>7</sup> Timothy A. Thornton,<sup>3</sup> Cathy C. Laurie,<sup>3</sup> Kenneth Rice,<sup>3</sup> and Xihong Lin<sup>1,\*</sup>

## ASSOCIATION STUDIES ARTICLE

### Genome-wide association study of iron traits and relation to diabetes in the Hispanic Community Health Study/Study of Latinos (HCHS/SOL): potential genomic intersection of iron and glucose regulation?

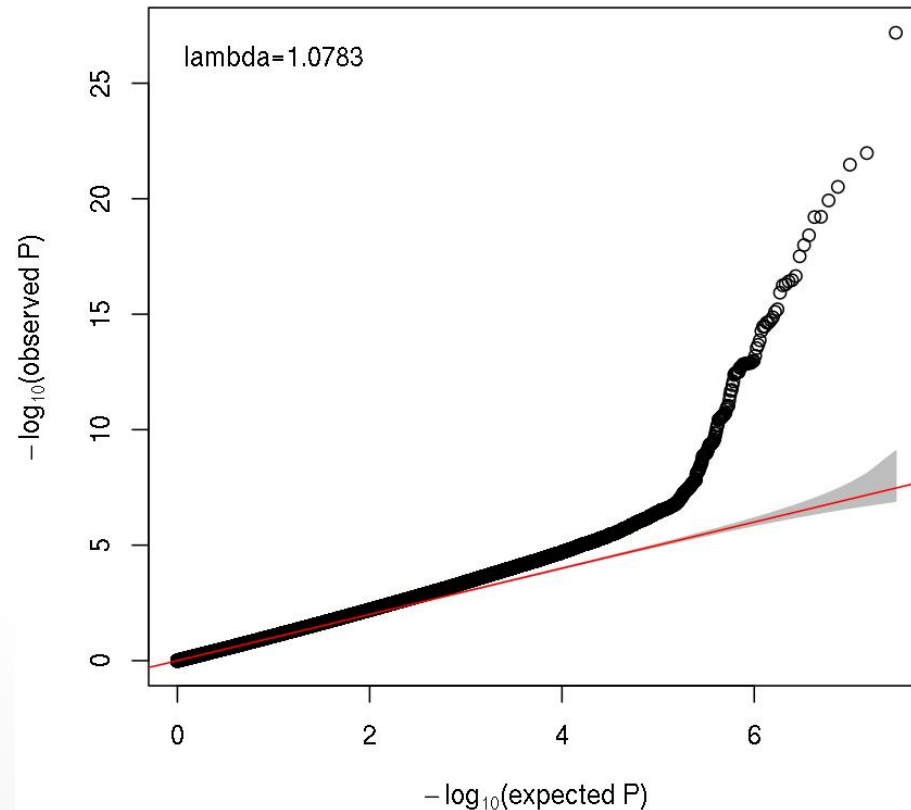
# Association Analysis with Linear Mixed Models

- LMM-OPS worked well for WHI-SHARe Hispanics and the Hispanic Community Health Study / Study of Latinos
- Applied LMM-OPS and standard LMM to PAGE and TOPMed traits
- Very perplexing results for many phenotypes!
- One phenotype with badly behaved results if fasting glucose

# PAGE Analysis: Standard LMM (EMMAX/GCTA)

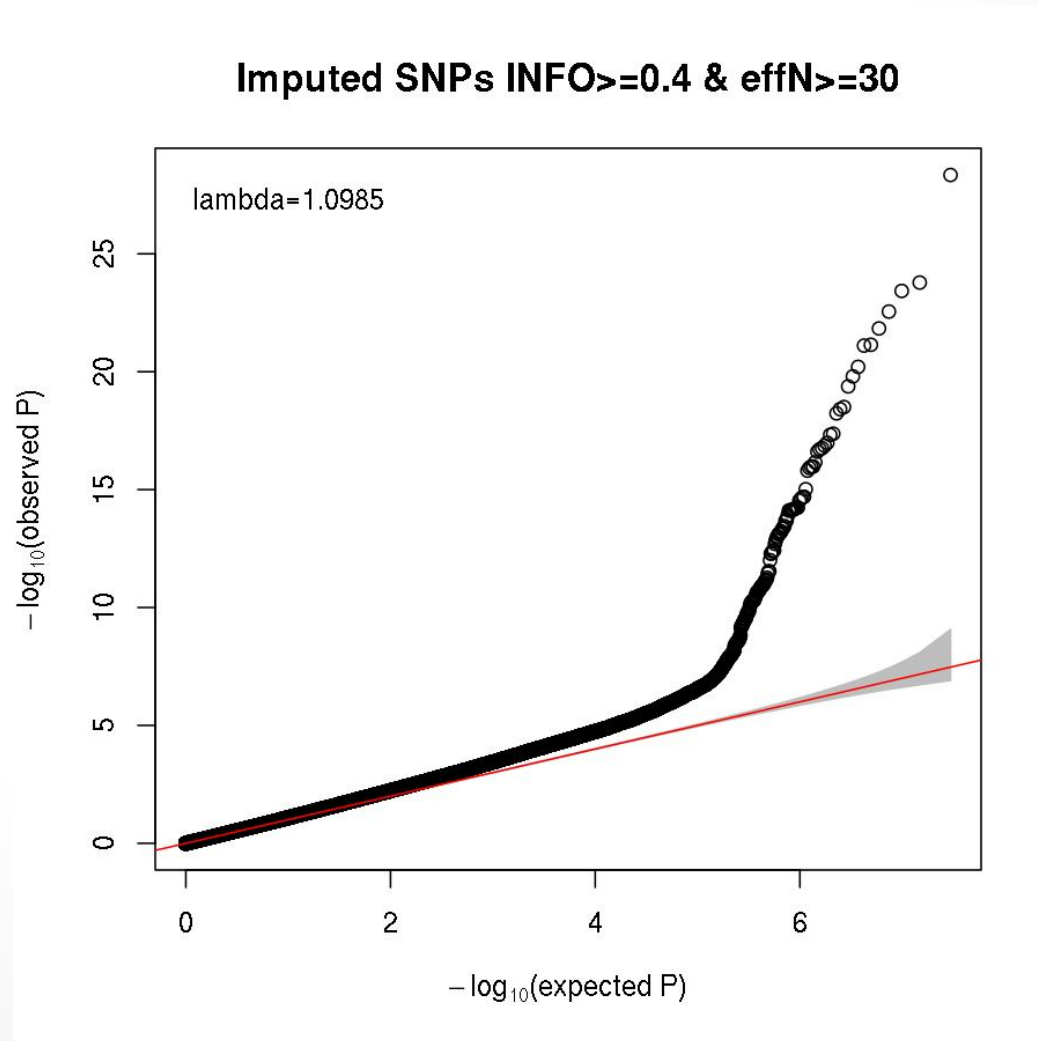
- Results for PAGE fasting glucose with adjustment for BMI

Imputed SNPs  $\text{INFO} \geq 0.4$  &  $\text{effN} \geq 30$



# PAGE Analysis: LMM-OPS

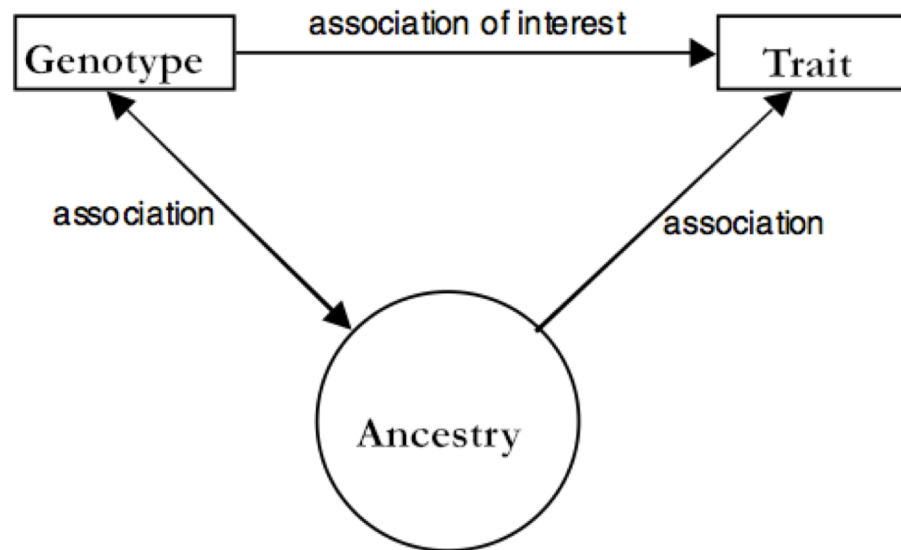
- LMM-OPS Results for fasting glucose with adjustment for BMI in PAGE





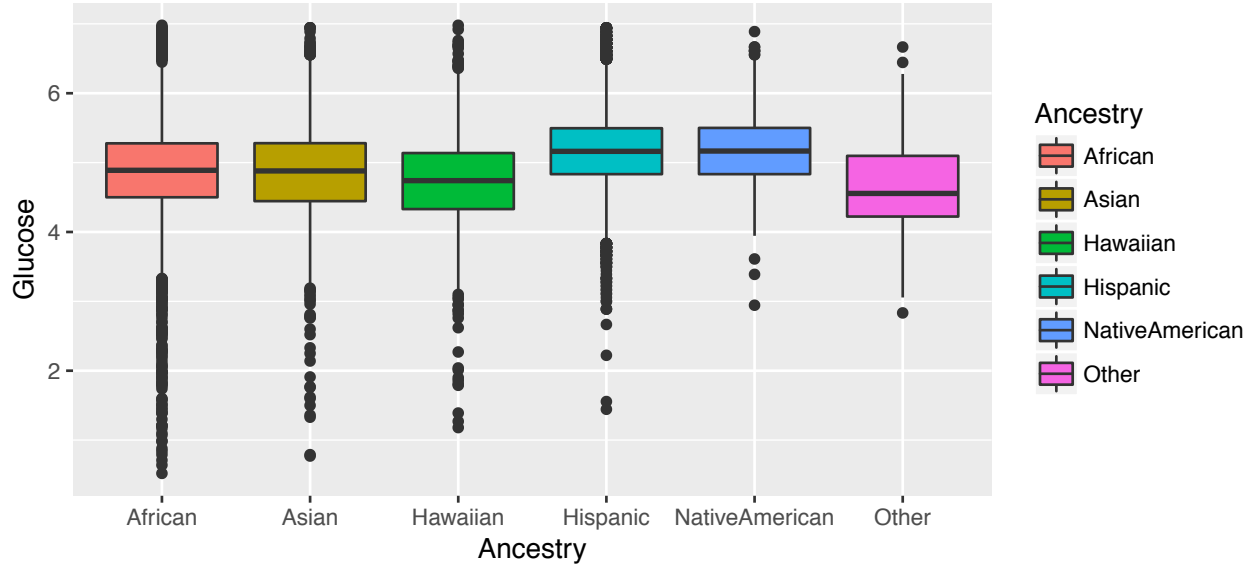
# Problems? Confounding due to combining samples from multi-ethnic populations

- Ethnic groups (and subgroups) often share distinct dietary habits and other lifestyle characteristics that result in traits of interest having **different distributions** that are correlated with genetic ancestry and/or ethnicity.

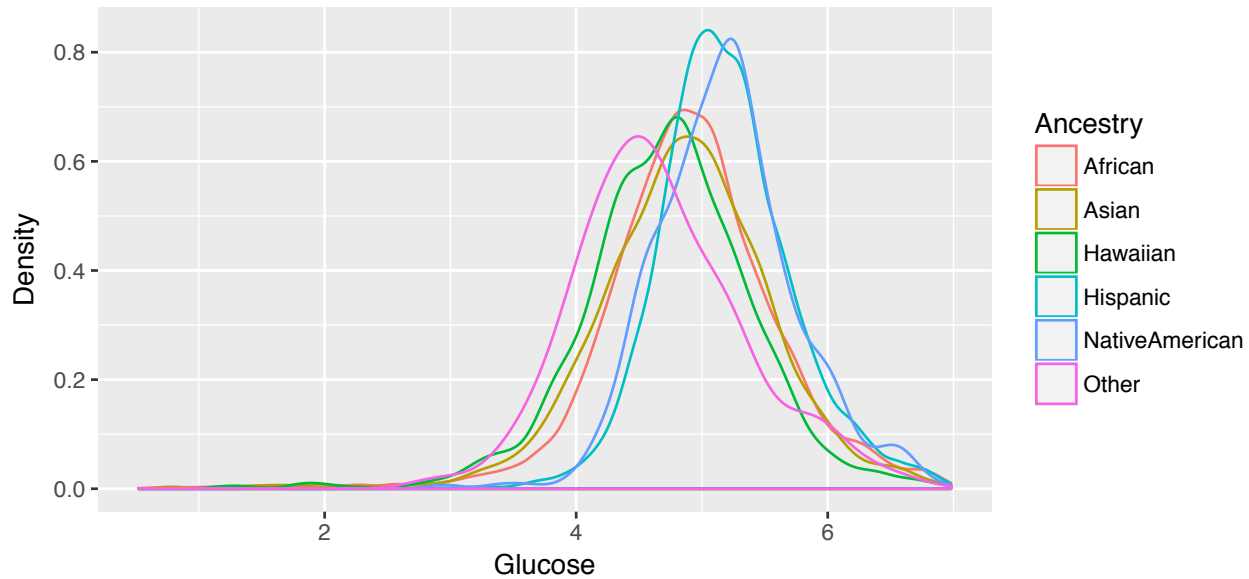


# PAGE Glucose by Race/Ethnicity

## PAGE Glucose Boxplots



## PAGE Glucose Densities



# PAGE Principal Components Analysis

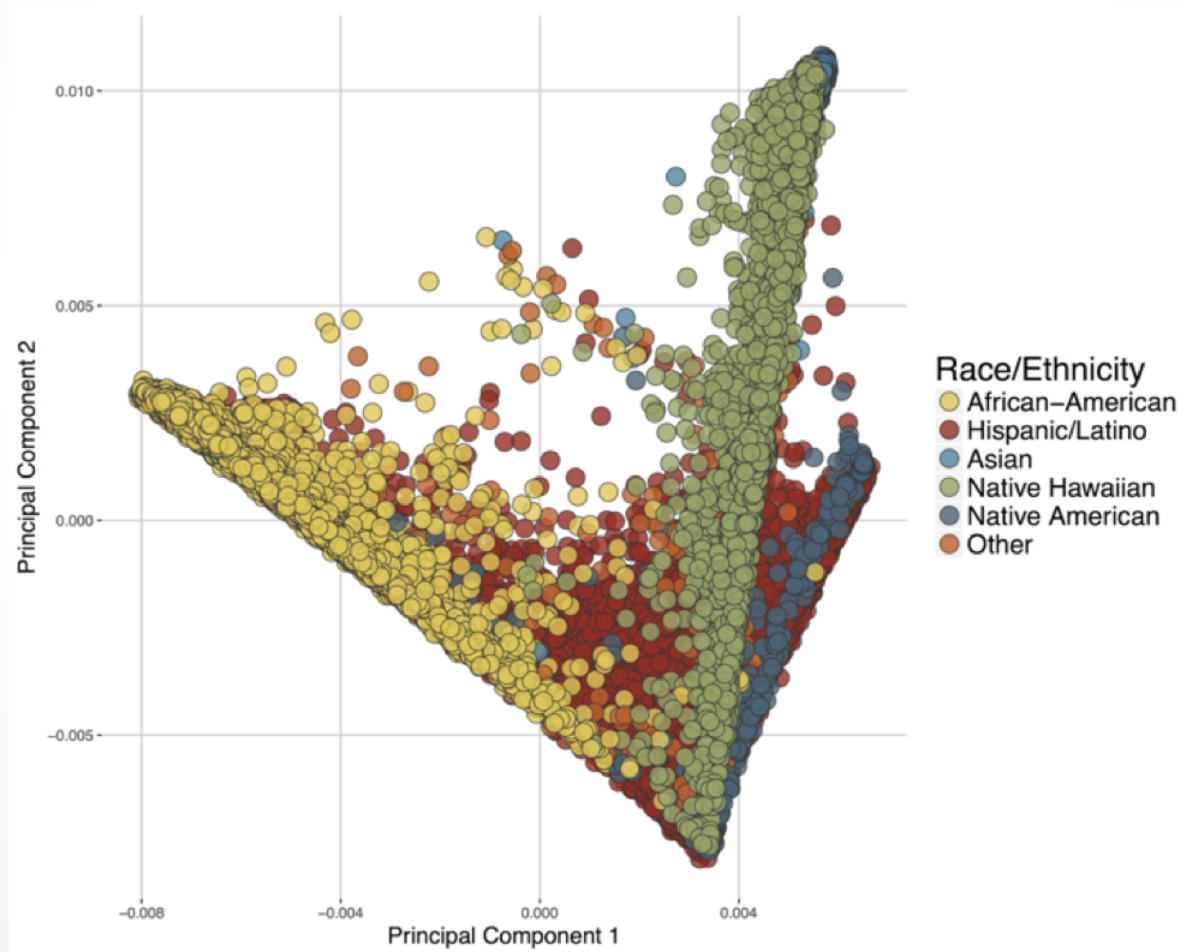


Figure Courtesy of Mariaelisa Graff  
(UNC Chapel Hill)

# PAGE Principal Components Analysis

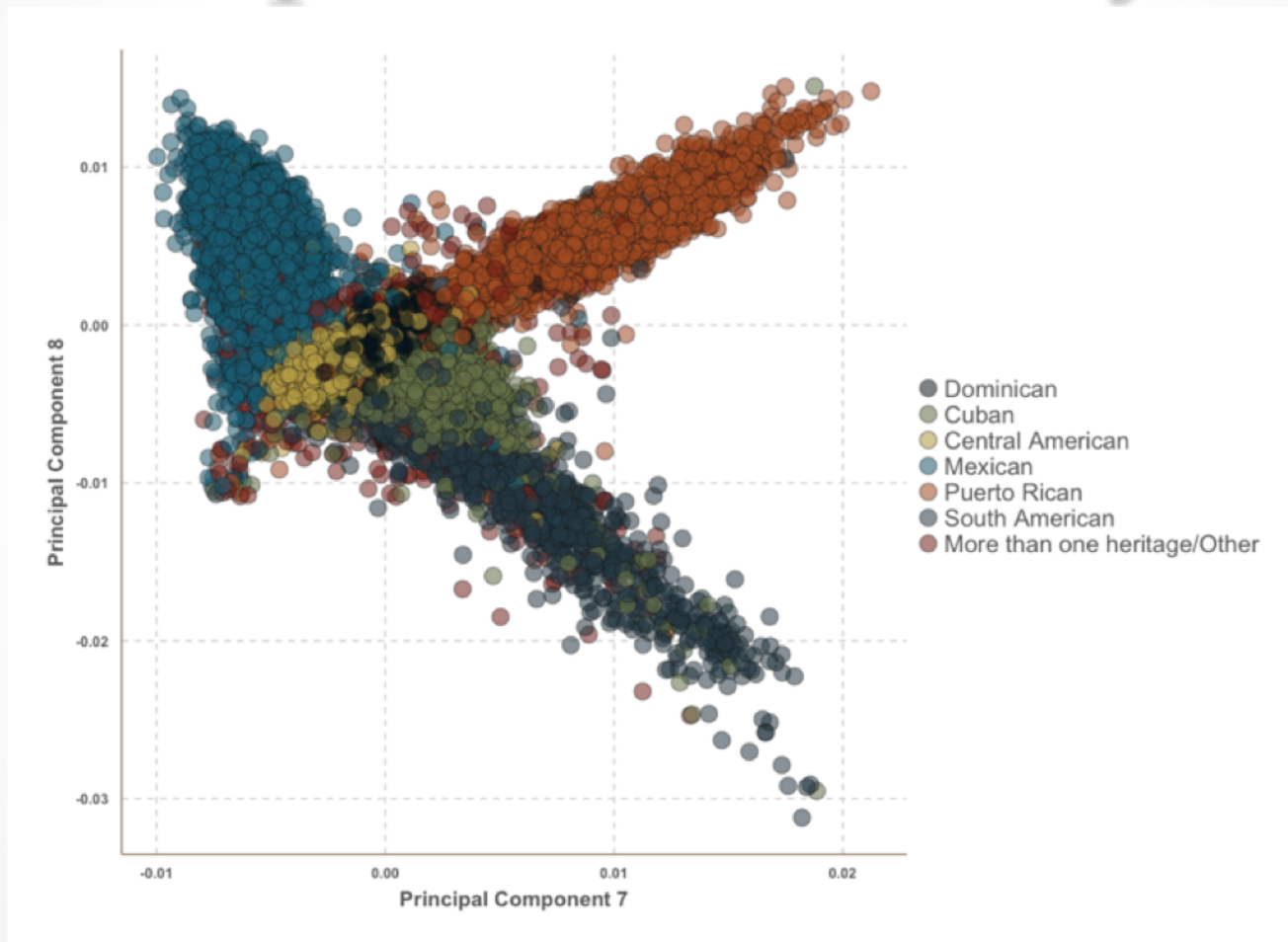


Figure Courtesy of Mariaelisa Graff  
(UNC Chapel Hill)

# PAGE Analysis: Standard LMM (EMMAX/GCTA)

- Recall the standard LMM model

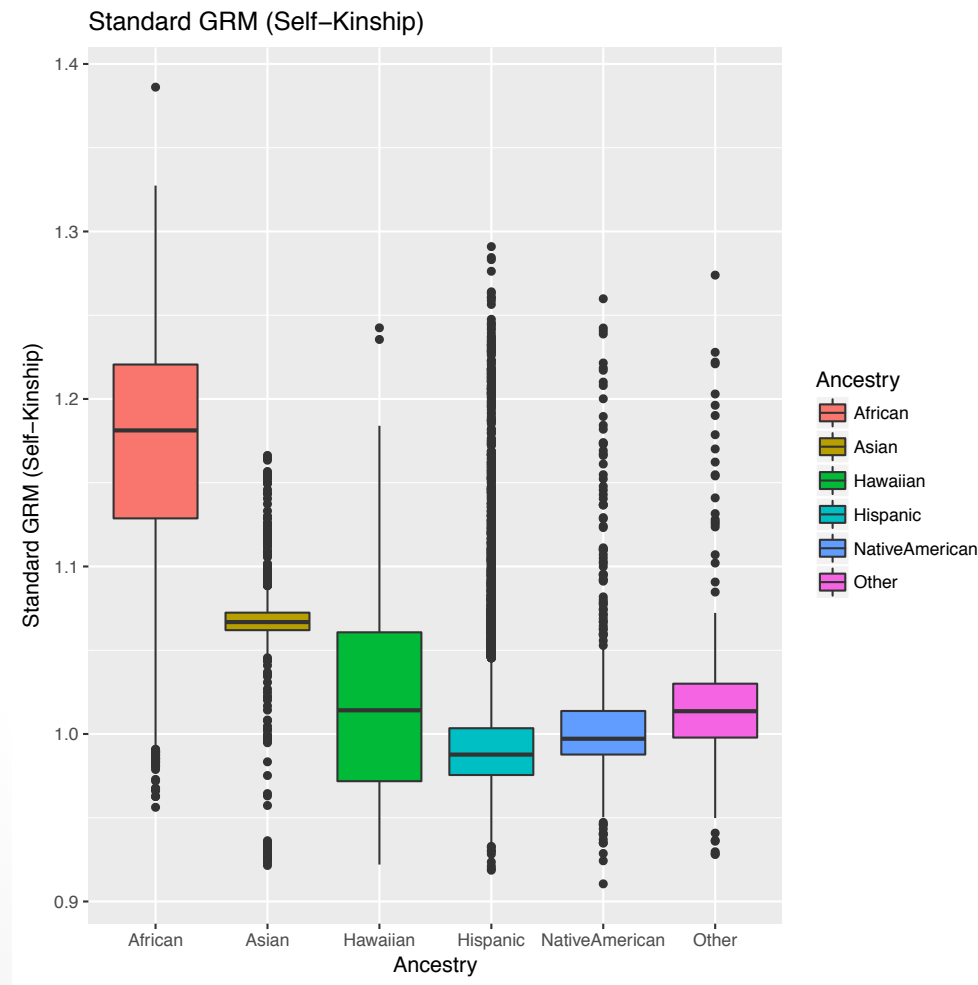
$$\mathbf{Y} = \mathbf{g}_s \beta_s + \mathbf{X} \boldsymbol{\alpha} + \boldsymbol{\epsilon} \quad \text{with} \quad \boldsymbol{\epsilon} \sim \text{N}(\mathbf{0}, \boldsymbol{\Sigma} \equiv \sigma_A^2 \boldsymbol{\Psi} + \sigma_\epsilon^2 \mathbf{I})$$

- For glucose with BMI adjustment, variance components were estimated
- Residual variance component:  $\hat{\sigma}_\epsilon^2 = 0.258$
- Additive genetic variance component:  $\hat{\sigma}_A^2 = 0.063$

$$\hat{\Psi}_{ij} = \frac{1}{M} \sum_{m=1}^M \frac{(g_{mi} - 2\hat{p}_m)(g_{mj} - 2\hat{p}_m)}{2\hat{p}_m(1 - \hat{p}_m)}$$

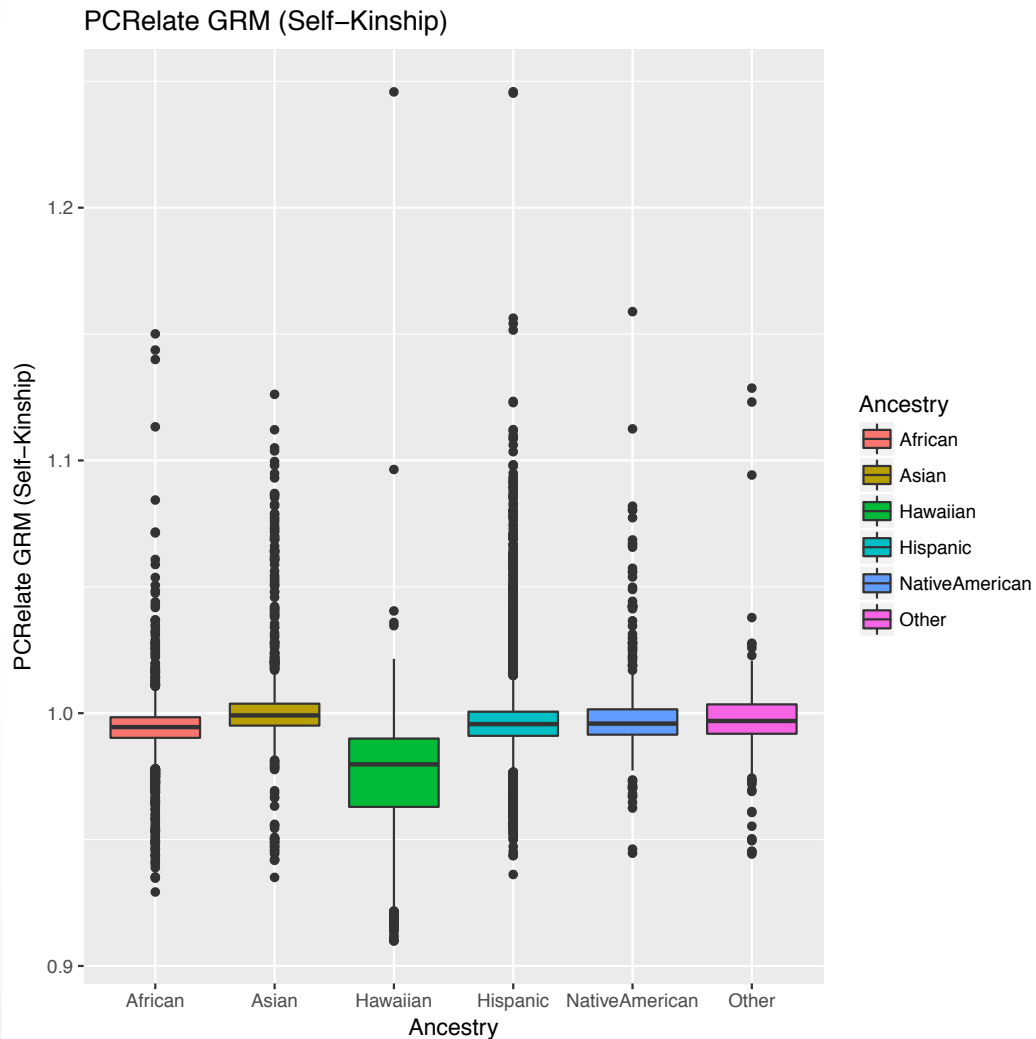
# Standard GRM for PAGE

- Below is a figure of the diagonal elements of the standard GRM for PAGE.



# LMM-OPS GRM for PAGE

$$Y = \mathbf{g}_s \beta_s + \mathbf{X} \alpha + \mathbf{V} \gamma + \epsilon \quad \text{with} \quad \epsilon \sim N(\mathbf{0}, \Sigma \equiv \sigma_A^2 \Phi + \sigma_\epsilon^2 \mathbf{I})$$



# Heterogeneity in Phenotypic Variances

- For multi-ethnic population samples, we should expect confounding due to traits having different means and variances.
- Developed **HEMMAT**: Heterogenous Effects Mixed Model Association Test.
- HEMMAT is an extension of LMM-OPS to allow for multiple random effects to be included in the model, in addition to the PC-adjusted GRM.
- HEMMAT incorporates additional **random effects to allow for heterogeneous variances** by self-reported race/ethnicity or by study (or any other discrete classification)



# Real Phenotype Data for Methods Comparison

- Variance components were estimated for BMI adjusted fasting glucose
- Additive genetic variance component is .032
- Residual variance components for each race/ethnicity:

Race/Ethnicity	Sample Size	Residual Variance
AA	6457	0.41
HS	13556	0.20
AS	1918	0.40
HI	1400	0.43
NA	412	0.24
Other	168	0.35

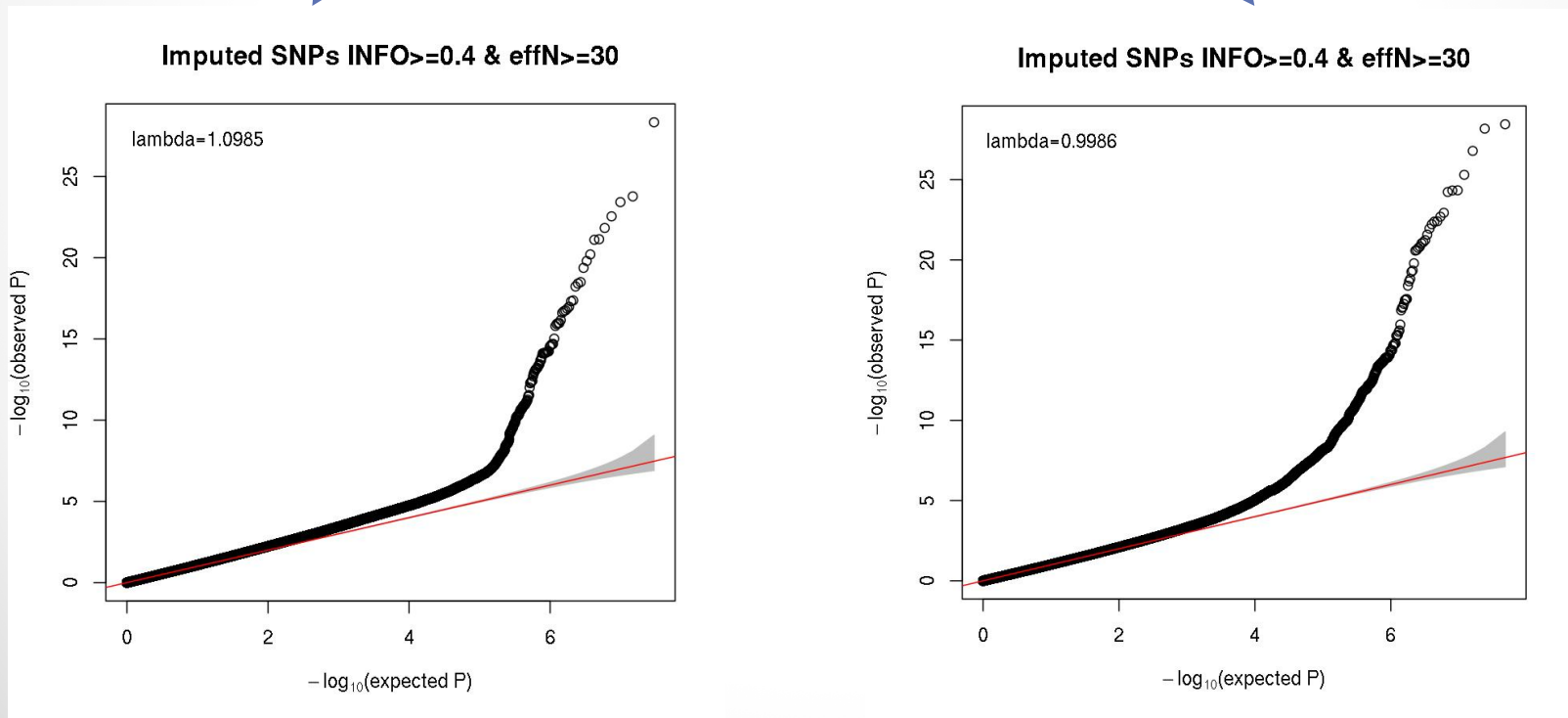
# HEMMAT Association Analysis Results

- HEMMAT results for fasting glucose with adjustment of BMI

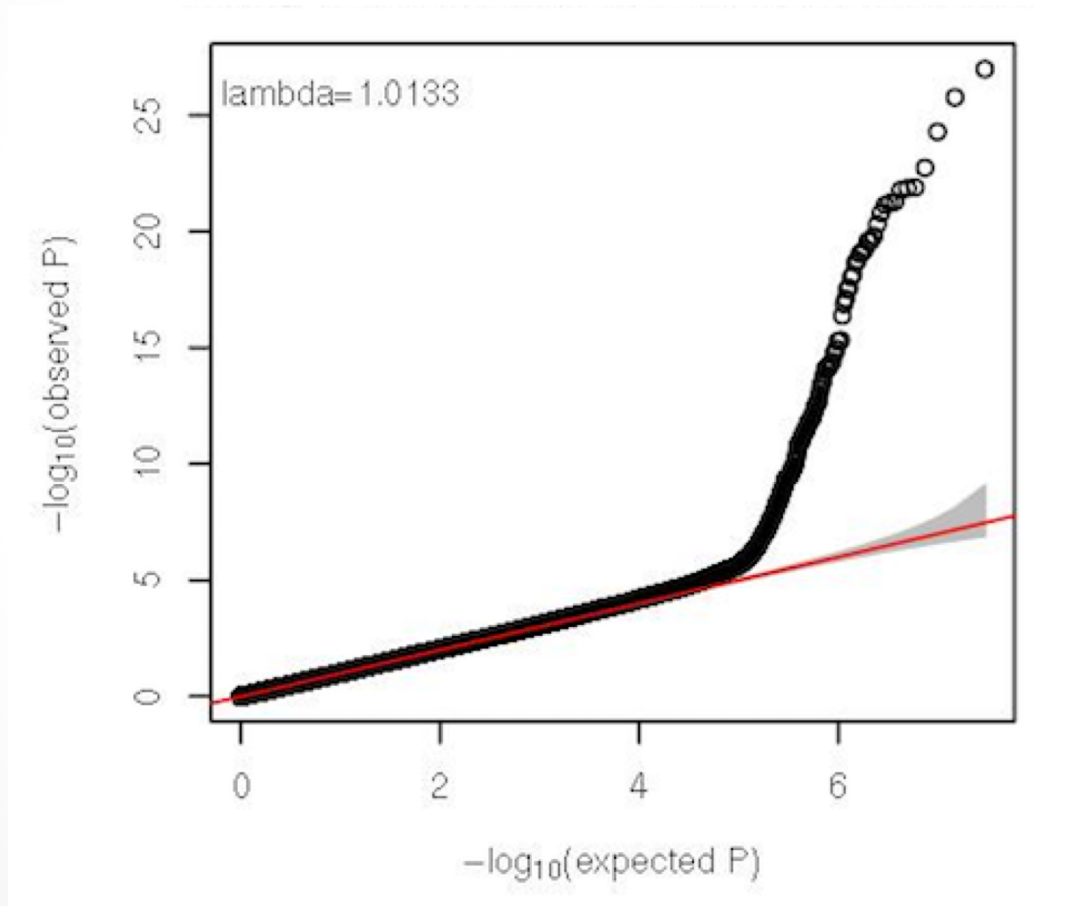
**Homogenous Trait Variance**



**Heterogeneous Trait Variance by Ethnicity/Race**



# Linear Mixed Model with Standard GRM and Heterogenous Variances



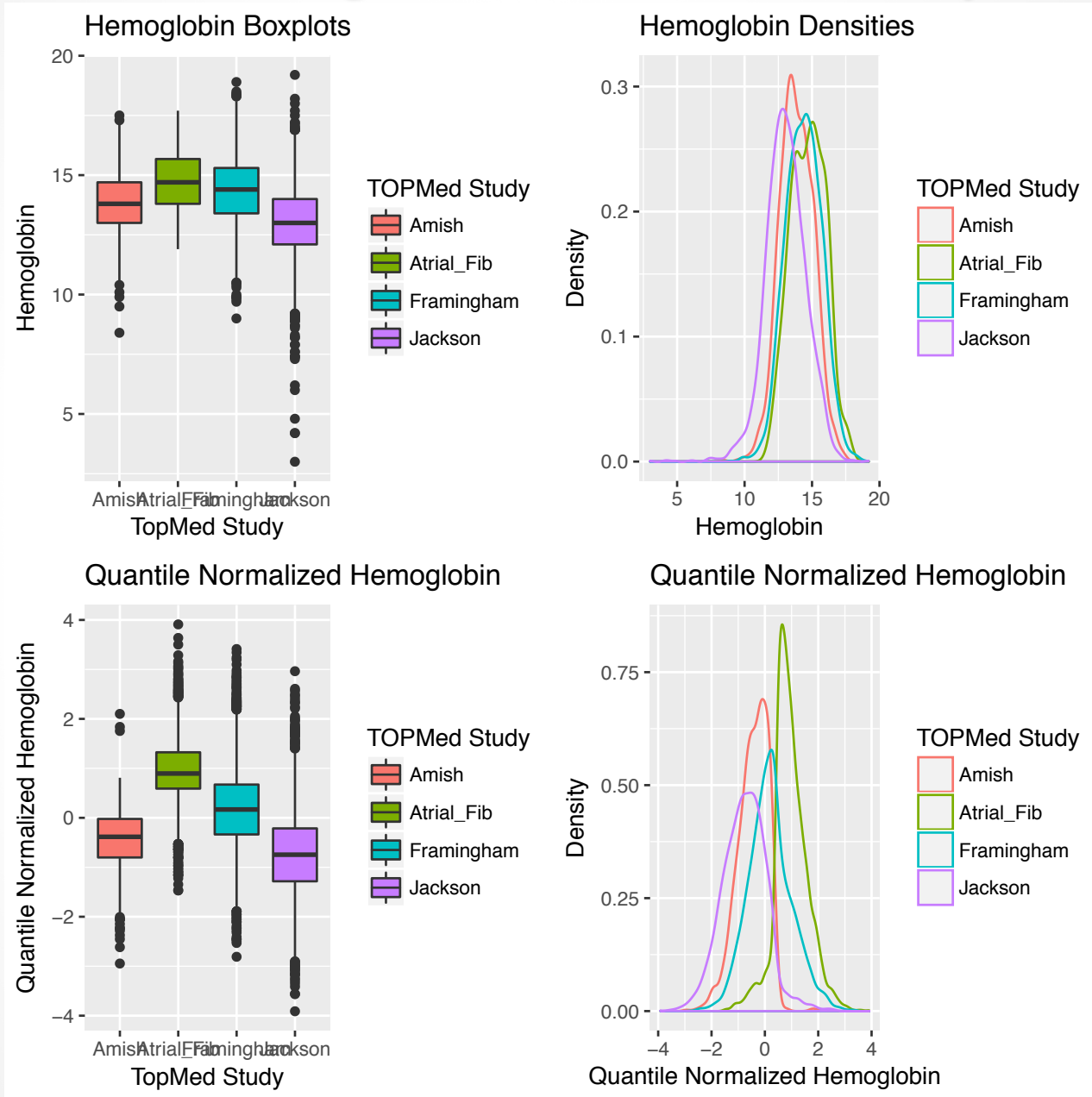
# Comparison of HEMMAT and Standard LMM with Heterogenous Variance

- Both HEMMAT and LMM with a standard GRM and heterogenous variances have proper control of genomic inflation for fasting glucose in PAGE
- HEMMAT: 168 genome-wide significant variants ( $p < 5e-08$ ) for fasting glucose identified
- LMM with standard GRM and heterogenous variance: 146 genome-wide significant variants for fasting glucose identified

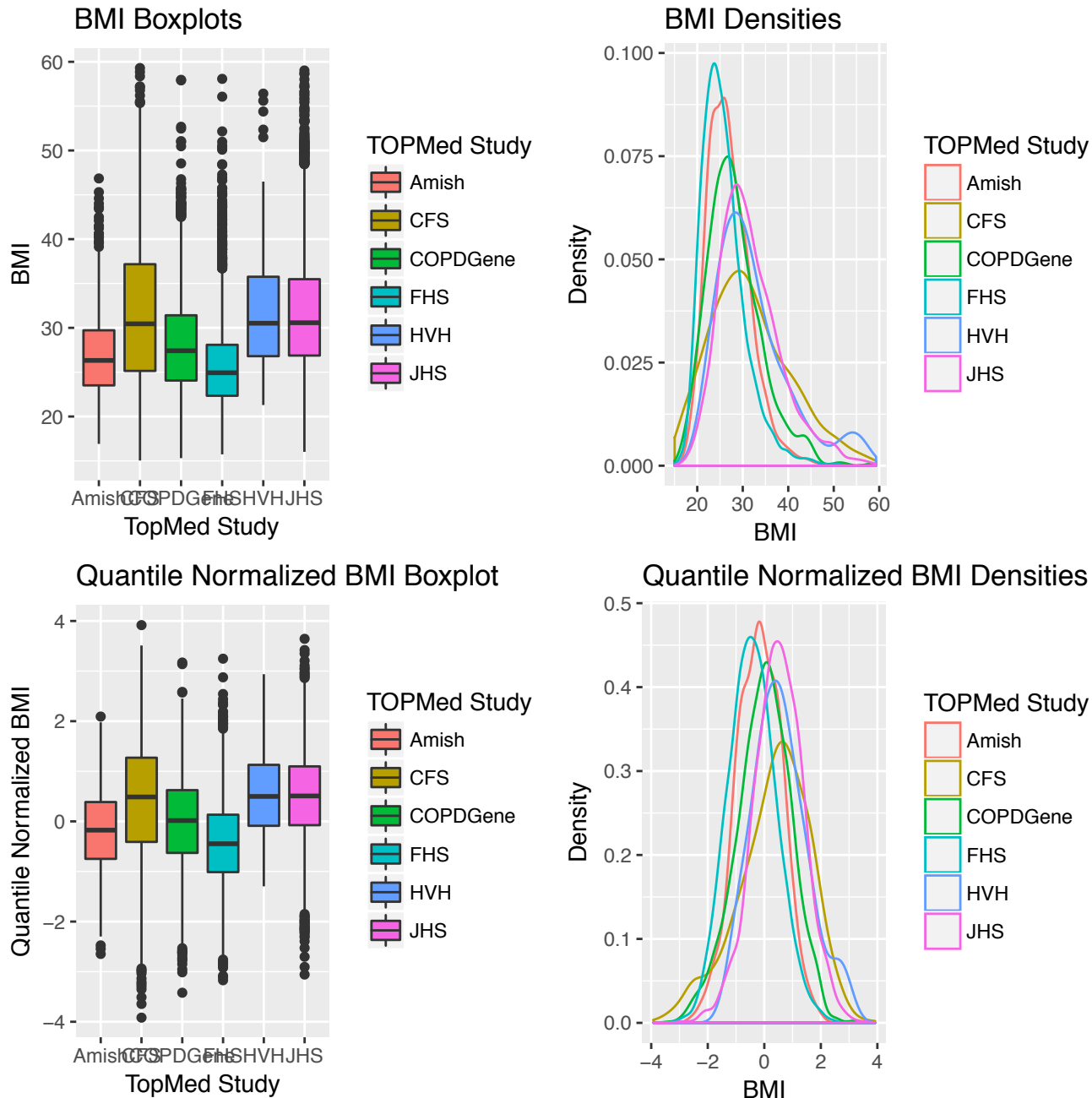
# Comparison of HEMMAT to LMM with Heterogenous Variance

LOCUS	Chr	Position	HEMMAT	LMM with Standard GRM and Heterogenous Variances
GCKR rs780093	2	27742603	1.03E-11	1.48e-10
ABCB11 rs563694	2	169774071	5.85e-25	1.64e-22
GCK rs1799884	7	44229068	3.94e-23	7.35e-22
TCF7L2 rs7903146	10	114758349	2.78e-08	1.47e-07
MTNR1B rs10830963	11	92708710	3.50e-29	4.48e-28
FOXA2 rs3833331	20	22562326	3.24e-12	3.83e-11

# TOPMed Hemoglobin Distributions by Study



# TOPMed BMI Distributions by Study



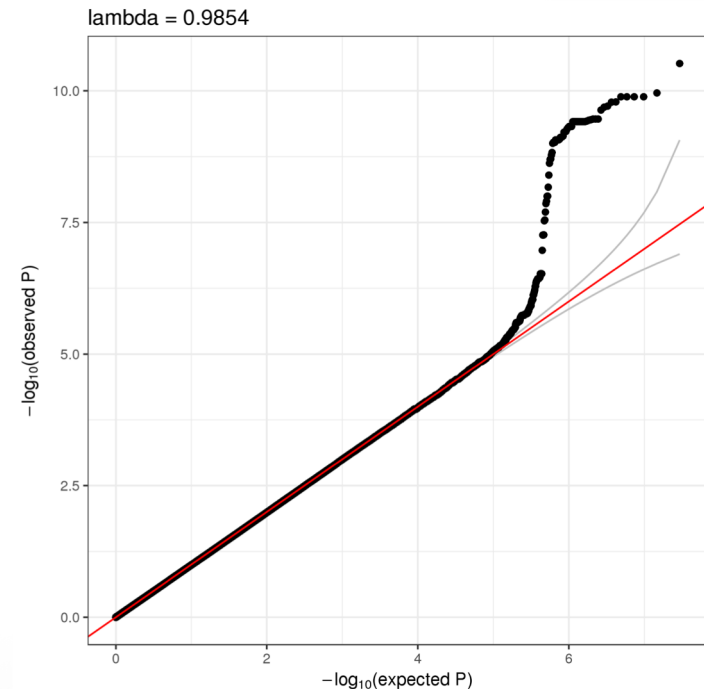
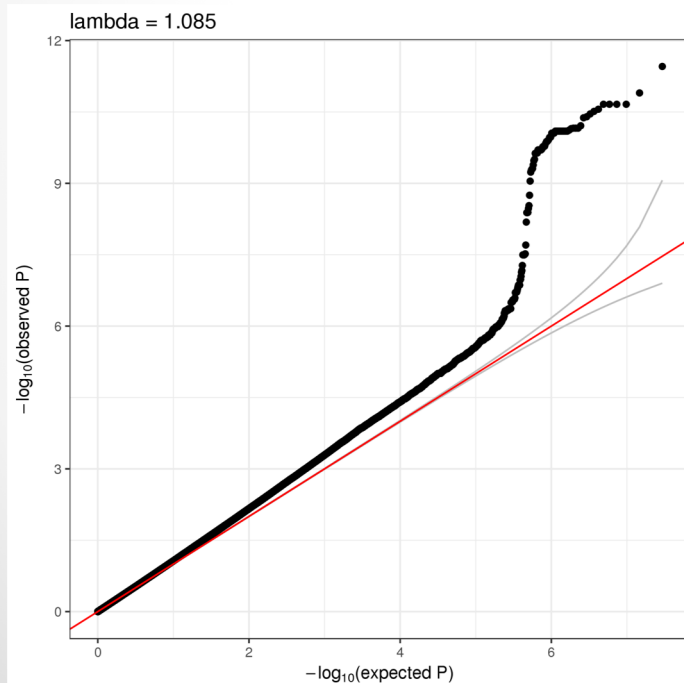
# Heterogeneity in Hemoglobin Phenotypic

## Variances: By Study

- HEMMAT results for hemoglobin: heterogeneous phenotypic variances

**Homogenous Residual Variance**

**Heterogeneous Residual Variances by Study**

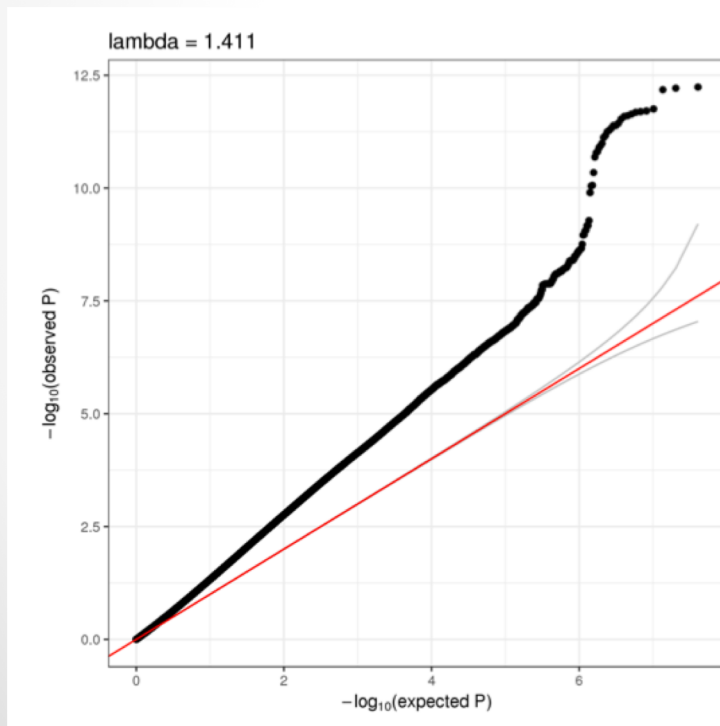




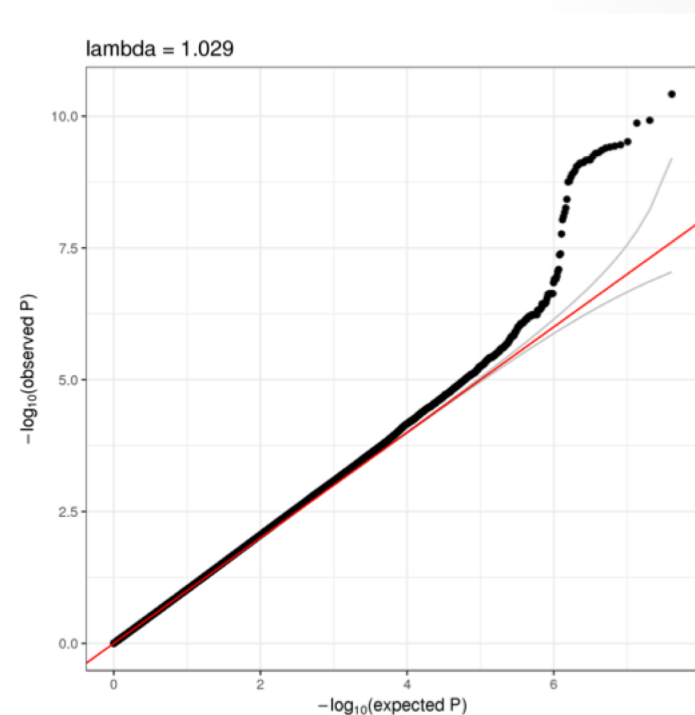
# Heterogeneity in BMI Phenotypic Variances

- HEMMAT Association results for BMI: heterogeneous phenotypic variances

**Homogenous Residual Variance**



**Heterogeneous Residual Variances by Study**



# Heterogenous residual variances for BMI

- Residual variance components of BMI for a few studies in TOPMed

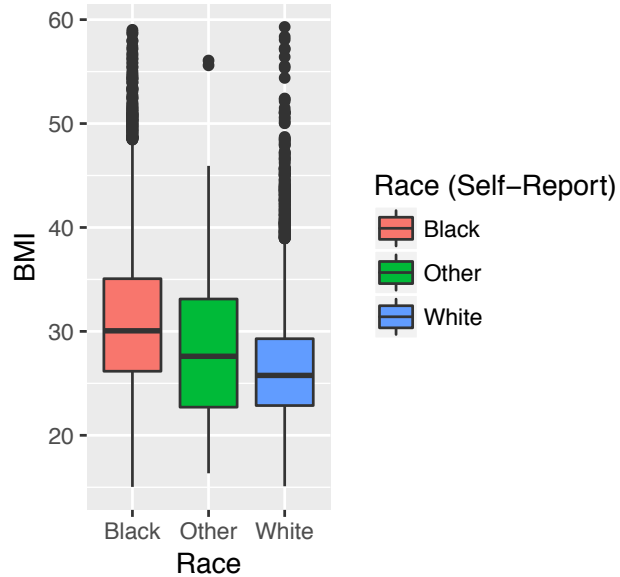
TopMed Cohort Study	Phenotypic Residual Variances
Jackson Heart Study	35.44
CFS	52.33
Framingham Heart Study	13.14
Amish	12.19
COPDGene	26.61
HVH	61.31

# Allowing for heterogeneity in variances: By Self-Reported Race

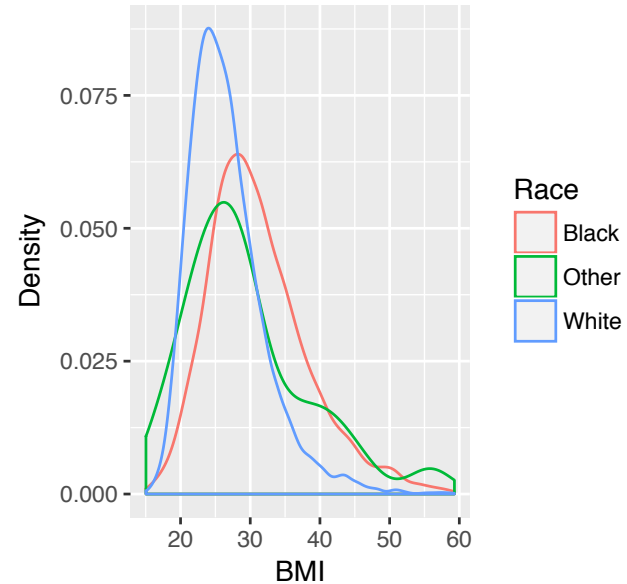
- There are limitations with modeling heterogeneous variances by study in TOPMed.
- A number of TOPMed studies have multiple ethnicities/ancestries.
- Also explored the differences in BMI distribution by self-reported race.

# TOPMed BMI Distributions: By Race

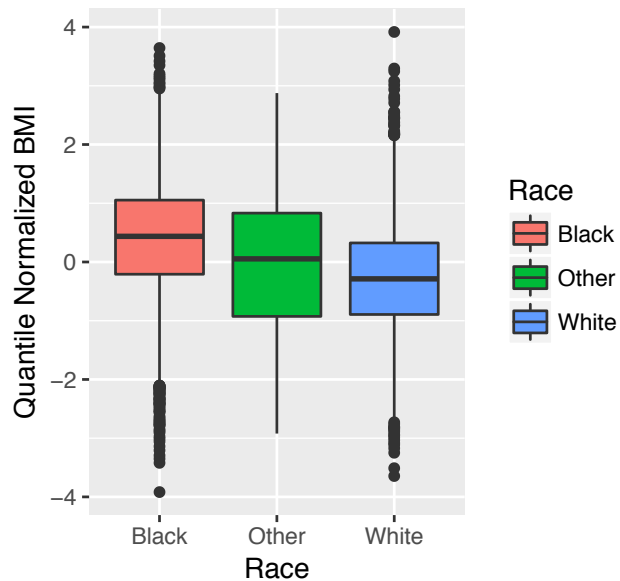
## BMI Boxplots by Self-Reported Race



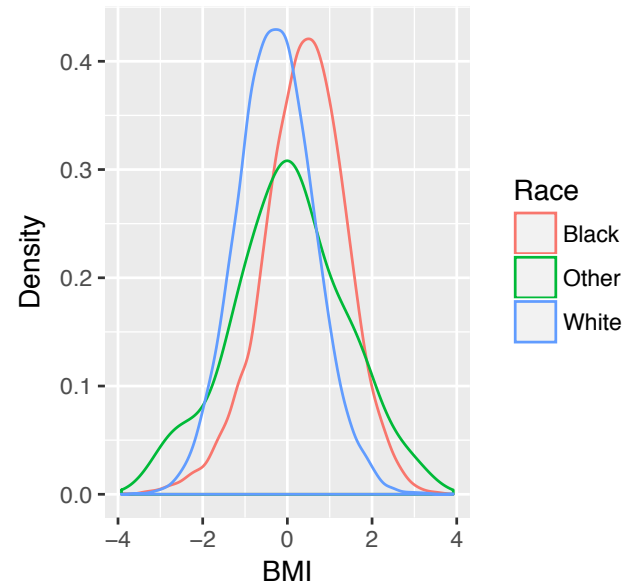
## BMI Densities



## Quantile Normalized BMI Boxplots



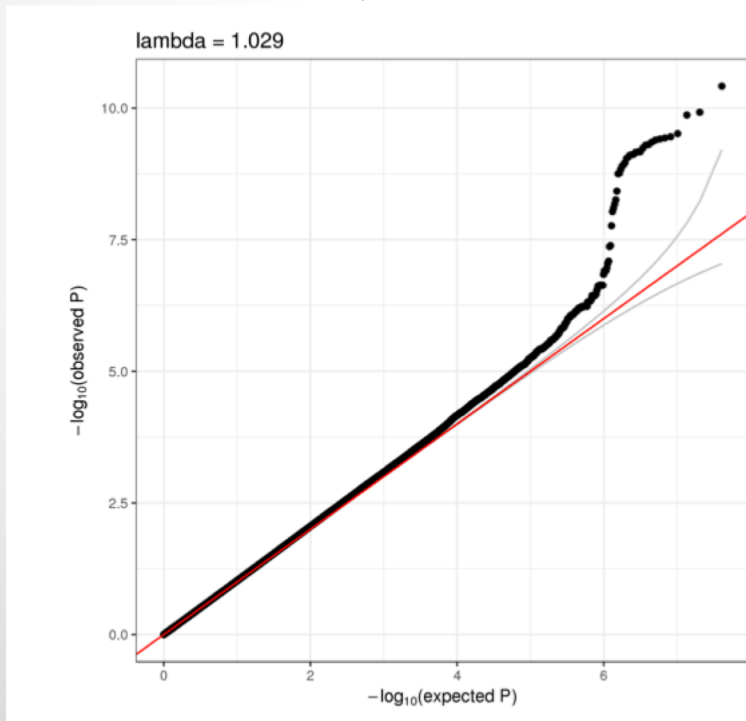
## Quantile Normalized BMI Densities



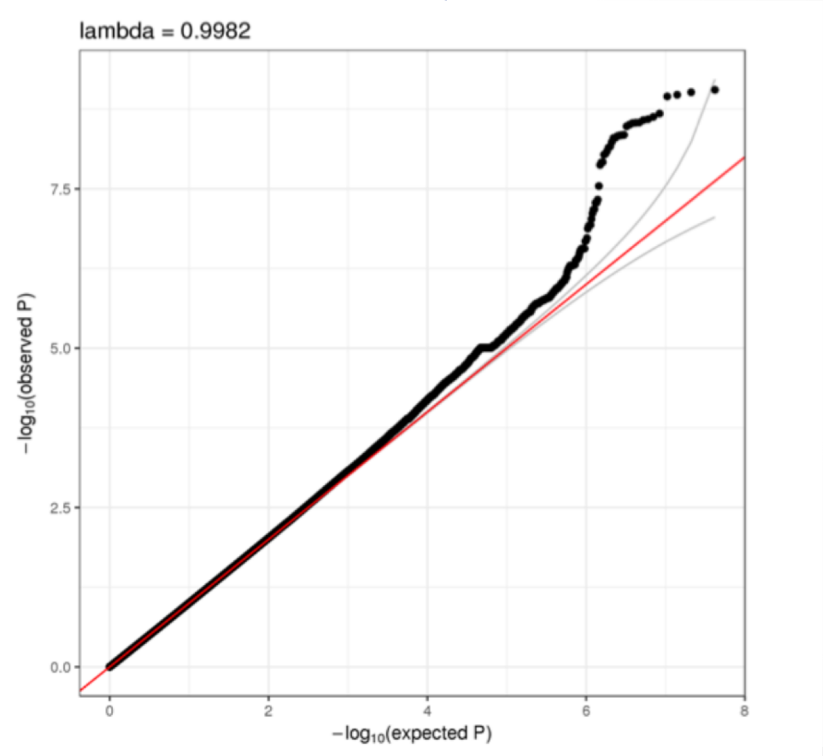
# BMI Heterogeneity: Study vs. Self-Reported Race

- HEMMAT Association results for BMI: heterogeneous phenotypic variances

**Heterogenous Residual Variance by Study**



**Heterogenous Residual Variance by Race**



# Novel Discoveries in Multi-Ethnic PAGE Study

- Conducted a GWAS of 26 clinical and behavioral phenotypes in 49,839 non-European individuals.
- Allowed for heterogeneous variances across all phenotypes
- 574 GWAS catalog variants across these traits were confirmed in PAGE
- 28 novel loci were identified.
- Manuscript submitted to *Nature*. Currently is in revision and will be re-submitted soon.



# GENESIS SOFTWARE

- **GENESIS:** R software package is available from Bioconductor
- Installation in R:
  - `source("https://bioconductor.org/biocLite.R")`
  - `biocLite("GENESIS")`
- Current release of GENESIS:
  - PC-AiR
  - PC-Relate
  - LMM-OPS
  - HEMMAT (for multi-ethnic populations with heterogenous variances)
  - Burden and SKAT tests have been extended, allow for population structure/relatedness and heterogenous variances

# Current/Future Work: X Chromosome analyses

HIGHLIGHTED ARTICLE  
GENETICS | INVESTIGATION

## Detecting Heterogeneity in Population Structure Across the Genome in Admixed Populations

Caitlin McHugh, Lisa Brown, and Timothy A. Thornton<sup>1</sup>

Department of Biostatistics, University of Washington, Seattle, Washington 98195

McHugh et al., *Genetics*, 2016





# Current/Future Work: X Chromosome analyses in TOPMed and PAGE

- Relatedness and population structure on the X-chromosome is quite different than the autosomes in multi-ethnic populations
- Dr. Caitlin McHugh has been developing extension of mixed methods for appropriate association testing (single variant and testing multiple variants simultaneously) on the X in multi-ethnic populations.
- Currently implementing the methods in GENESIS and will be applying to TOPMed and PAGE.



# Acknowledgements: TOPMed

## UW DCC for TOPMed

Department of Biostatistics, University of Washington

Bruce Weir

Tim Thornton

Ken Rice

Sharon Browning

Brian Browning

Katie Kerr

Tamar Sofer

Adam Szpiro

Cathy Laurie

David Levine

Cecelia Laurie

Stephanie Gogarten

Adrienne Stilp

Deepti Jain

Quenna Wong

Xiuwen Zheng



National Heart, Lung,  
and Blood Institute

**TOPMed**

## Grant Funding:

NIGMS: P01 GM099568

NHLBI: R01 HL120393,  
HHSN268201300005C

# Acknowledgements: PAGE

- PAGE Members & PAGE Lead Coordinators.
- The Population Architecture Using Genomics and Epidemiology (PAGE) program is funded by the National Human Genome Research Institute (NHGRI) and the National Institute on Minority Health and Health Disparities (NIMHD), supported by U01HG007417 (BioME), U01HG007416 (CALiCo), U01 HG007397 (MEC), U01 HG007376 (WHI), and U01HG007419 (Coordinating Center).

# Special Acknowledgements

- Matthew Conomos, PhD; UW GAC



- Stephanie Gogarten, PhD; UW GAC



- Xiuwen Zheng, PhD; UW GAC

